

UNCLASSIFIED

AD NUMBER
AD407434
NEW LIMITATION CHANGE
TO Approved for public release, distribution unlimited
FROM Distribution authorized to U.S. Gov't. agencies and their contractors; Administrative/Operational Use; MAY 1963. Other requests shall be referred to Army Electronics Command, Army Liaison Group [Project MICHIGAN], The University of Michigan, P.O. Box 618, Ann Arbor, MI.
AUTHORITY
USAEC ltr, 10 Nov 1966

THIS PAGE IS UNCLASSIFIED

2900-281-T

CATALOGED BY DDC

407 434

AS AD No.

407 434



Report of Project MICHIGAN

CONTINUOUS HUMAN ESTIMATION OF A TIME-VARYING, SEQUENTIALLY DISPLAYED PROBABILITY

GORDON H. ROBINSON

RECEIVED
MAY 15 1963
TISIA B

ENGINEERING PSYCHOLOGY LABORATORY

Institute of Science and Technology

THE UNIVERSITY OF MICHIGAN

May 1963

Contract DA-36-039 SC-78801

NO. OTS

2900-281-T

Report of Project MICHIGAN

**CONTINUOUS HUMAN ESTIMATION OF A
TIME-VARYING, SEQUENTIALLY
DISPLAYED PROBABILITY**

GORDON H. ROBINSON

May 1963

Engineering Psychology Laboratory
Institute of Science and Technology
THE UNIVERSITY OF MICHIGAN
Ann Arbor, Michigan

NOTICES

Sponsorship. The work reported herein was conducted by the Institute of Science and Technology for the U. S. Army Electronics Command under Project MICHIGAN, Contract DA-36-039 SC-78801. Contracts and grants to The University of Michigan for the support of sponsored research by the Institute of Science and Technology are administered through the Office of the Vice-President for Research.

Note. The views expressed herein are those of Project MICHIGAN and have not been approved by the Department of the Army.

Distribution. Initial distribution is indicated at the end of this document. Distribution control of Project MICHIGAN documents has been delegated by the U. S. Army Electronics Command to the office named below. Please address correspondence concerning distribution of reports to:

Commanding Officer
U. S. Army Liaison Group
Project MICHIGAN
The University of Michigan
P. O. Box 618
Ann Arbor, Michigan

ASTIA Availability. Qualified requesters may obtain copies of this document from:

Armed Services Technical Information Agency
Arlington Hall Station
Arlington 12, Virginia

Final Disposition. After this document has served its purpose, it may be destroyed. Please do not return it to the Institute of Science and Technology.

PREFACE

Project MICHIGAN is a continuing, long-range research and development program for advancing the Army's combat-surveillance and target-acquisition capabilities. The program is carried out by a full-time Institute of Science and Technology staff of specialists in the fields of physics, engineering, mathematics, and psychology, by members of the teaching faculty, by graduate students, and by other research groups and laboratories of The University of Michigan.

The emphasis of the Project is upon research in imaging radar, MTI radar, infrared, radio location, image processing, and special investigations. Particular attention is given to all-weather, long-range, high-resolution sensory and location techniques.

Project MICHIGAN was established by the U. S. Army Signal Corps at The University of Michigan in 1953 and has received continuing support from the U. S. Army. The Project constitutes a major portion of the diversified program of research conducted by the Institute of Science and Technology in order to make available to government and industry the resources of The University of Michigan and to broaden the educational opportunities for students in the scientific and engineering disciplines.

Progress and results described in reports are continually reassessed by Project MICHIGAN. Comments and suggestions from readers are invited.

Robert L. Hess
Director
Project MICHIGAN

ACKNOWLEDGMENTS

The author is grateful to Professor Elmer G. Gilbert for his guidance and instruction during the author's graduate program, and to Professor Paul M. Fitts for his considerate leadership.

Professor Arthur W. Melton, Director of the Engineering Psychology Laboratory, provided support and direction for the laboratory in which this work was performed.

Special gratitude must be extended to Professor Ward Edwards for his direct supervision of the work and his many suggestions concerning interesting directions of investigation. The author also greatly appreciates the intellectual stimulation and leadership that Professor Edwards has provided during the past three years.

CONTENTS

Notices	ii
Preface	iii
Acknowledgments	iv
List of Figures	vi
List of Tables	vii
Glossary	viii
Abstract	1
1. Introduction	1
2. The Experiment	4
2.1. The Task	4
2.2. Input Selection	7
2.3. Flash Series Generation	7
2.4. Experimental Variables	8
2.5. Task Instructions	9
3. Experimental Results	10
3.1. The Pilot Experiment	10
3.2. Response Measures	11
3.3. Data Analysis	13
3.4. Experimental Data	13
3.5. Summary of Results	24
4. Mathematical Models	25
4.1. A Model With Geometric Weighting	26
4.2. A Model With Constant Weighting	32
4.3. A Descriptive Model	35
4.4. Normative Variation	40
5. Discussion	42
6. Conclusion	46
Appendix A: Input Probability Generation	47
Appendix B: Variances of Sample Averages from Finite Populations	48
Appendix C: Order of Presentation	50
Appendix D: Instructions	51
Appendix E: Two Qualitative Response Exceptions	52
Appendix F: Data Not Averaged over Subjects	52
References	60
Bibliography	60
Distribution List	62

FIGURES

1. Tracking Console	5
2. Analog Computer Circuit for Payoff	9
3. A Typical Response to a Subproblem	10
4. Detection as a Function of Step Size and Flash Rate	14
5. Percentage of Subproblems in Which "No Detection" Occurs, as a Function of Step Size and Flash Rate	14
6. Detection as a Function of Probability, Constraint, and Small- and Large-Step Problems	15
7. Detection as a Function of Step Size, Sample Rate, and Constraint	15
8. Convergence as a Function of Step Size and Flash Rate	17
9. Percentage of Subproblems in Which "No Convergence" Occurs, as a Function of Step Size and Flash Rate	17
10. Percentage of Subproblems in Which "Accidental Initial Convergence" Occurs, as a Function of Step Size and Flash Rate	17
11. Convergence as a Function of Probability, Constraint, and Small- and Large-Step Problems	17
12. Convergence as a Function of Step Size, Flash Rate, and Constraint	18
13. Root Mean Square Error over the Whole Subproblem as a Function of Flash Rate	19
14. Mean Error as a Function of Probability	20
15. Root Mean Square Error as a Function of Step Size and Constraint	21
16. Root Mean Square Error as a Function of Flash Rate and Constraint	21
17. Root Mean Square Error as a Function of Probability and Constraint	22
18. False Alarm Rate in False Alarms Per Flash as a Function of Step Size and Constraint	22
19. False Alarm Rate in False Alarms Per Flash as a Function of Flash Rate and Constraint	23
20. False Alarm Rate in False Alarms per Flash as a Function of Probability and Constraint	23
21. Responses of Two Mathematical Models and a Subject to a Portion of a Large-Step Problem, Random Constraint, at 1 FPS	41
22. Responses of Two Mathematical Models and a Subject to a Portion of a Small-Step Problem, Random Constraint, at 1 FPS	42
23. Detection as a Function of Step Size for Four Subjects	53
24. Detection as a Function of Flash Rate for Four Subjects	54
25. Convergence as a Function of Step Size for Four Subjects	55
26. Convergence as a Function of Flash Rate for Four Subjects	56
27. Root Mean Square Error as a Function of Flash Rate for Four Subjects	57

28. Root Mean Square Error after Convergence, as a Function of Flash Rate for Four Subjects	58
29. False Alarm Rate, in False Alarms Per Flash, as a Function of Flash Rate for Four Subjects	59

TABLES

I. Parameter Sets for the Descriptive Model Yielding Minimum Values of $\langle e_{ms}^2 \rangle$	38
II. Comparison of the Performances of the Subjects and the Mathematical Models.	45
III. Probabilities Used to Generate Sequences of Flashes	47

GLOSSARY

a	a constant
C	convergence, a measure of the subject's response (see page 11)
\bar{C}	no convergence, a measure of the subject's response (see page 11)
D	detection, a measure of the subject's response (see page 11)
\bar{D}	no detection, a measure of the subject's response (see page 11)
$e(n)$	a model's error at n, $e(n) = r(n) - P_2$
$\overline{e^2(n)}$	the expected value of $e^2(n)$, an ensemble average
$\langle e^2 \rangle_A$	the average value of $e^2(n)$ over some set of samples A
$e_M(n)$	the descriptive model's error at n, $e_M(n) = r(n) - P_2$
$e_{MS}(n)$	the error (i.e., disagreement) between the subject and the descriptive model at n, $e_{MS}(n) = R_n - r(n)$
$e_S(n)$	the subject error at n, $e_S(n) = R_n - P$
$\langle e_S^2 \rangle_P, \langle e_M^2 \rangle_P, \langle e_M e_{MS} \rangle_P$	the average values of $e_S^2(n)$, $e_M^2(n)$ and $e_M(n)e_{MS}(n)$ over a problem
E(x)	the expected value of x
FAR	false alarm rate, a measure of the subject's response (see page 12)
fps	flashes per second, the flash rate
IC	initial convergence, a measure of the subject's response (see page 11)
k	$\frac{1}{T} \sum_{i=1}^{i=g} (P_{i-1} - P_1)^2$
k_1	the decision criterion level in the descriptive model
k_2	the fractional response adjustment in the descriptive model
k_3	the number of flashes in u(n) in the descriptive model
k_4	the flash shift between the subject's and the descriptive model's responses
M	the last flash in a subproblem
M_i	the length of subproblem i
ME_C	mean error after convergence, a measure made on the subject's response (see page 12)
n	a flash index, n = 1 is the first flash in a subproblem
N	the total number of flashes in the model's summation
P	the probability of a 1 in a 0,1 binary series
P_1	the probability of a 1 in the binary series i
Pr(x)	the probability of x
r	the geometric ratio in the geometric model

$r(n)$	a model's response at n
$\overline{r(n)}$	the expected value of $r(n)$, an ensemble average
R_n	the subject's response at n
RMSE	the square root of the mean square error, a measure of the subject's response (see page 11)
RMSE _C	the square root of the mean square error after the point of convergence, a measure of the subject's response (see page 12)
ρ	a correlation coefficient
S	the total number of subproblems in a problem
s_n	the n -th flash in a binary series
σ_1^2	the variance of a sample from a population with $P = P_1$
$\sigma_r^2(n)$	the variance of $r(n)$ at n , an ensemble average
$\sigma_s^2(n)$	the variance of s_n at n , an ensemble average
T	the total number of flashes in a problem
$u(n)$	an average of k_3 flashes in the descriptive model
v	$\frac{1}{S} \sum_{i=1}^{i=S} \sigma_i^2$
w_1	a weight attached to sample s_{n-1+1}

CONTINUOUS HUMAN ESTIMATION OF A TIME-VARYING, SEQUENTIALLY DISPLAYED PROBABILITY

ABSTRACT

This experiment examines the human ability to give a direct magnitude estimate of a time-varying probability. The subject positioned a "tracking" lever at his estimate of the current mean of a sequentially displayed binary distribution. The distribution samples were presented at a fixed rate by two flashing lights. The distribution mean changed in step increments of varying size and spacing. The experimental variables included flash rate and a constraint on the randomness of the flash series.

Detailed measures were made of both the transient and static responses to each step change. The transient response was more rapid and consistent than had been anticipated and occurred with step changes as small as 0.12. The average static response showed no systematic bias as a function of probability and had an RMS error approximately equal to that of a 17-sample average.

Two simple mathematical models are derived to provide quantitative comparisons with the subjects' data. A descriptive model is also derived which satisfies some basic properties of the task behavior. The parameters for this model are selected for two specific experimental situations.

1

INTRODUCTION

Human decision tasks can be described as static or dynamic. In a dynamic decision task, some of the relevant stimuli vary as a function of time, or of past decisions, or of both. The decision maker must keep track of these changes in order to perform satisfactorily.

This experiment examines the human ability to follow or estimate a time-varying probability, which could be an important input to a dynamic decision task. The experiment attempts to isolate the estimation of the probability from the use of the estimate in making decisions. The task selected was the estimation of the mean of a binary (Bernoulli) distribution. Samples (0 or 1) from the distribution were displayed sequentially, and the subject continuously estimated the distribution mean, which varied with time. The experiment is described in detail in Section 2.

The study of probability estimation isolated from decision making is important for two reasons. First, in decision making under uncertainty, the estimation of probabilities is always at least an implicit part of the task. A decision maker's ability to produce decisions which maximize expected value will depend directly on his ability to estimate the probabilities of the various alternative courses of action.

Second, there is an applied interest in the human ability to turn uncertainty into probability. In many systems involving stochastic inputs, it is relatively easy to automate the application of decision rules. It is far harder, however, to find automatic means for supplying the probabilities and payoffs necessary for the application of the rules. Probability estimation, then, is a candidate for inclusion as a human task in semiautomatic information processing and decision-making systems in which the subsequent choice of a course of action is performed automatically.

RESEARCH ON ESTIMATION AND PREDICTION

Human binary-choice behavior has been studied extensively. This experiment complements previous studies by isolating the estimation function and by using changing probabilities.

Most studies of human binary choice do not include estimation as an explicit part of the subjects' task. These studies usually generate prediction data, which are then averaged over blocks of trials (decisions, choices) to produce prediction frequencies. A prediction frequency of 0.67 on trials 121 through 150 would indicate that the predictions during these 30 trials were distributed about two to one between the two choices. The subject may or may not be told the correct choice after each trial. The correct choices are drawn in some manner from a stationary binary distribution.

Examples of these experiments, often called probability learning experiments, are reported by Grant, 1953 [1]; Hake and Hyman, 1953 [2]; Hake, 1954 [3]; Estes, 1957 [4]; and Neimark and Shuford, 1959 [5]. Most of these studies report prediction frequencies asymptotically approaching the frequency of correct choice or the generating probability. This phenomenon has been named "probability matching." This behavior is not optimum. The optimum strategy, under instructions to maximize correct choices, is to predict consistently the more probable event. This event can be inferred from the relative frequency of previous events.

Behavior significantly different from matching has been reported by Gardner, 1959 [3], and Edwards, 1961 [7]. The number of trials may have been insufficient in some of the experiments in which matching was found. An unpublished experiment by Tannenbaum and Edwards at The University of Michigan indicates that the amount of reward for a correct choice interacts with the prediction frequency. Some subjects used near-optimum strategies.

A few studies have looked at the estimation ability of the binary decision maker. Grant [1] reports an experiment by Hornseth in which the subject was asked to guess, at the end of 150 choice trials, which event had been the more frequent. The prediction frequencies for the last 30-trial block were close to the matching level. The data on guessing the overall frequency were plotted as the percentage of correct guesses. These data showed that the percentage of

correct guesses at a particular event frequency was higher than the event frequency (i.e., an event frequency of 0.70 would be guessed to be the more frequent over 70 percent of the time).

Grant concludes from this experiment and presumably from other probability learning experiments that the processes of estimation and prediction are distinct, and that prediction is the more accurate. The processes may indeed be distinct, but the accuracy of estimation was not measured by Hornseth's experiment.

Hake [3] surveys a major portion of the probability learning literature (including experiments by Estes, Grant, Hornseth, Hake, and Hyman), and concludes that estimation is not accurate enough to be the basis for binary predictions.

Subjects in these experiments should have based their choices on estimates of the event frequencies or generating probabilities. To conclude that the non-optimum performance was an indicator of inaccurate probability estimates is unjustified, however.

Neimark and Shuford [5] included estimation as an explicit part of the task in a probability learning experiment. Besides making a choice at each trial, some subjects were required to estimate the proportions of the past events. The event frequency was 0.67. These subjects gave unbiased estimates and had prediction frequencies significantly higher than the matching level, whereas subjects who only predicted produced frequencies at the matching level. These results suggest that explicit estimation improved prediction.

Erlick [8] looked at estimation without a decision task. He presented 100 binary events at a rate of five per second and asked for an estimate of the more frequent event and for an actual estimate of the event frequency on a continuous scale. Four event frequencies were used: 0.50-0.50; 0.48-0.52; 0.45-0.55; and 0.43-0.57. The data indicated that the more frequent event was selected correctly 75% of the time when the event-frequency difference was approximately 0.08 (0.46-0.54). For 0.50 and 0.52, the median estimate of the frequency was within 0.01; for 0.55 and 0.57 the median estimate was about 0.02 high.

All of the experiments reviewed above used stationary processes to generate the binary events. A few experiments have used a dynamic generating process, but since prediction was the required task in all of these, they give only indirect evidence on estimation.

Grant [1] reports an experiment in which the generating probability changed periodically as a square wave. The probability values always differed by 0.50 with higher values: 1.00, 0.90, 0.80, and 0.70. The period was 40 events, and two and one half cycles were presented. A prediction frequency was calculated by averaging over five trials and about 40 subjects. This prediction frequency followed the cyclic change only when the higher probability was 1.00 or 0.90, and reached 0.95 in 20 trials at 1.00, and 0.70 in 20 trials at 0.90. Apparently no systematic

performance changes occurred during the two and one half cycles. The subjects were evidently not instructed about the nonstationarity of the generating process. Such instructions would probably have had an appreciable effect.

Goodnow and Pettigrew [9] performed a binary prediction experiment in which a change from 0.50-0.50 to 0-1.00 occurred. They found that the response to such a change was more rapid when the subjects had initially experienced a 0-1.00 series prior to the 0.50-0.50 series. Again, however, no specific instructions were given concerning the nonstationarity of the generating process.

In both Grant's and Goodnow's experiments there is evidence that a change in the generating probability of 0.50 will produce an appropriate change in the prediction frequency if the change is to an extreme probability. These extreme probabilities (1.00 and 0.90) evidently represent changes which are obvious even when the instructions induce no expectation of change.

Flood [10] discusses the strategy of a subject who may not be convinced that the probabilities are stationary. No particular results were obtained in an experiment designed to induce certainty versus uncertainty in the stationarity of a stationary generating process.

The human ability to estimate directly the magnitude of a stationary binary probability is uncertain. Most experimenters have postulated estimation only as an intervening variable between the display and a decision task. Decision behavior was improved in one experiment by including explicit estimation in the task. Two questions seem appropriate: what role does estimation play in a decision task, and how well can this estimation be performed? The experiment reported here sheds light on the second question as well as providing a fairly comprehensive look at the continuous estimation of dynamic probabilities.

2

THE EXPERIMENT

2.1. THE TASK

The task studied in this experiment was to estimate the mean of a binary distribution as samples (individual drawings) from that distribution were sequentially displayed. This task was selected for two reasons. First, it is completely described by one parameter, its mean. It is thus easily understood by people unfamiliar with the mathematical aspects of probability. Second, it can be readily related to the literature on binary decision and estimation, discussed in Section 1.

The display and response mechanisms, shown in Figure 1, were designed for convenient and effective control and interpretation. As the subject sat at the apparatus, samples from a

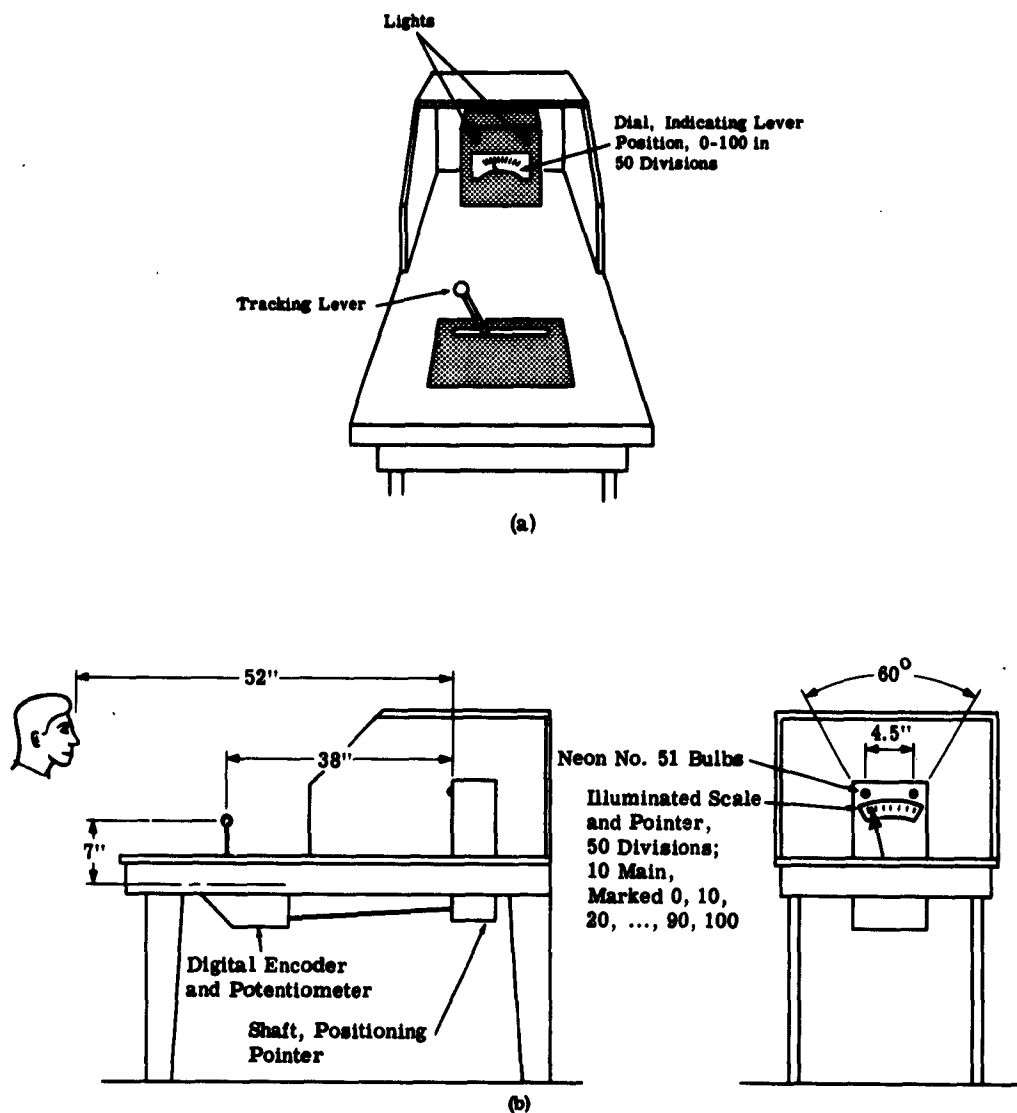


FIGURE 1. TRACKING CONSOLE. (a) Sketch. (b) Schematic. The flash had a duration of approximately 0.020 seconds. The intensity was adjusted to provide a clear indicator without glare. The room illumination was low.

binary distribution were presented to him at a fixed rate by two flashing lights. He reported his estimate of the mean of the distribution by moving the tracking lever so that the illuminated dial beneath the lights would indicate it. Thus the apparatus was a continuous-response mechanism—appropriate for the estimation of a continuously varying stimulus.

The lever was free to move between stops at 0 and 100 on the scale. The smallest scale division was 2. The variable error of the pointer was about one half of the least division, corresponding to a probability change of 0.01. A main scale division occurred at every fifth small division and was marked 0, 10, 20, . . . , 90, 100. The lever and associated mechanisms contained enough Coulomb friction to retain a setting without constant force; neither springs nor viscous friction was used.

The position of the lever was recorded by two means, Friden punched paper tape and a Sanborn continuous-strip recorder. The paper tape was punched in a Grey, or cyclically permuted binary code using six channels of an eight-channel punch to encode 101 symbols, 0 through 100. Pilot studies indicated that the rate of output sampling necessary to recover the response information depended on the flash rate, and that a response sampling rate equal to the flash rate would be adequate. Thus a sample was taken every two seconds at the slowest presentation rate, 0.5 flash per second, and every 0.125 second at the fastest presentation rate, 8 flashes per second. The punched-tape record was later transferred to IBM cards on a modified IBM tape-to-card converter, and the data analyzed on an IBM 709 data processing system. The Sanborn records were used in making qualitative judgments about the response and in selecting appropriate criteria for the computer analysis. They also permitted continuous monitoring of the task as it was performed.

The subject and his console were isolated in a small room. The subject wore noise-insulating ear muffs. He had a two-way communication system with the experimenter. A low-level white noise was presented by the earphones during the experimental run. When the experimenter spoke to the subject, the noise was automatically switched off. The subject's microphone was always on, and comments during the experimental run were permitted.¹

The task has a strong resemblance to a standard unidimensional manual tracking task. The main difference is the presentation of the target: instead of being displayed explicitly as a dot or a line, it exists only as a parametric description of the method used to select the flash sequence. In this experiment the generating process was time-variant, and the target could be defined as the mean of the distribution from which the last flash was drawn. It is impossible to recover the precise target from the information available to the subjects. The cursor, or 0-100

¹ Few comments were made; most of these were not printable.

dial pointer in this case, is pointed at an estimated value of the target. The system is essentially open-loop, since the lack of an explicit target prohibits the formation of an error signal. The dynamics are almost entirely in the mental computation; there was no indication that motor skill was a limiting factor.

The use of a tracking lever as a response means is unique in research on probability estimation. It is appropriate to the task and permits an easy understanding of the response scale by the subjects. Both end points are well fixed, in the same sense that impossible events and sure ones are fixed in value on a personal or subjective probability scale. The 50 point on the scale might also be considered as an anchor point, since all subjects clearly understood that 50% meant equally frequent flashes.

2.2. INPUT SELECTION

The input probability changed in a series of discrete steps. This input form permitted visual, qualitative interpretations to be made from the data in addition to the more extensive analysis done by the computer. The input form also permitted static as well as dynamic measurements to be made. The step-change sizes and their directions, as well as the number of flashes between steps, were selected randomly from a finite set of values. The sequence of steps so generated is called a problem. The mechanism for the generation of the sequences is described in detail in Appendix A.

Preliminary investigations revealed that step changes ranging from 0.06 to 0.64 in eight values would adequately cover the interesting range of probability change.² The number of flashes between step changes was selected from a set ranging from 34 to 89 flashes; the smallest number of flashes required to minimize interaction between successive step changes was 34. The range between 34 and 89 was considered sufficient to prevent any performance improvement due to the learning of inter-step length. A step change and the flashes until the next change are called a subproblem.

2.3. FLASH SERIES GENERATION

The flashes were drawn from finite populations without replacements. The population size was an experimental variable and is discussed below. Finite populations were selected to fix the average value of the flashes for each subproblem. The effects of finite population sampling on variances are shown in Appendix B.

²A pilot experiment with a simplified apparatus was run before the main console was built in order to establish the general form of the response and reasonable ranges for the independent variables. It is discussed in more detail in Section 3.

2.4. EXPERIMENTAL VARIABLES

Five independent variables were used in the experiment: the rate at which the flashes were presented; the magnitude and sign of each step change; the probability after the step change; a constraint on the randomness of the flash series; and subjects. The number of flashes between step changes was not studied as a variable. The variables had the following values:

Rate: 0.5, 1.0, 2.0, 4.0, and 8.0 flashes per second (fps)

Step size: 0.06, 0.12, 0.16, 0.18, 0.24, 0.32, 0.48, 0.64 (both + and -)

Probability: 0.02, 0.08, 0.14, 0.18, 0.26, 0.32, 0.34, 0.44, 0.50, and the complementary values between 0.50 and 1.00

The step changes and probability values were arranged in two problem types, described in detail in Appendix A. For one type, the small-step problems, the mean step change is approximately 0.15; for the other, the large-step problems, the mean step change is approximately 0.40. Both types contain the entire range of probabilities and are symmetric about 0.50.

The constraint variable had two values, leading to the random and the constrained problem types. The random problems were generated from finite populations which had the length of the respective subproblems being generated. These finite populations were thus of size 35 through 89 flashes. These sizes were considered large enough to yield experimental results fairly close to those which would result from infinite populations.

The constrained problems were generated from finite populations of 17 flashes. The lengths of the subproblems were arranged in whole-number multiples of 17: 34, 51, 78, and 85 flashes. It was assumed that this constraint would be sufficient to indicate those aspects of performance that constraint would affect. It is not a severe enough constraint to be readily perceived from inspection of the flash series, however. The same series of steps and probabilities were used in the random and in the constrained problems.

Each of the four subjects performed the task in 15 sessions, and saw the same series of problems in the same order. Each session lasting for about an hour, consisted of two or three problems separated by a short rest period.

Rates, small- and large-step problems, subjects, and constraints were exhaustively combined. The order of presentation was chosen at random under the constraint that the tracking sessions were of about the same length. (Appendix C gives the sequence used.) The pilot experiment had indicated that about 25 minutes, at two flashes per second, was the maximum time that a subject could be expected to track without a significant decrement in his performance. The problems presented at 0.5 and 1.0 flashes per second were given in four and two

separate sessions, respectively, in order to limit all sessions to a maximum of 25 minutes of continuous tracking.

2.4. TASK INSTRUCTIONS

Careful attention was paid to the instruction of the subjects prior to the recorded experimental sessions. This effort was repaid by an excellent consistency in the tracking behavior of the eight subjects, four in the pilot and four in the main experiment. (The standard instructions used are shown in Appendix D.) These served only as the initial formal introduction, however. Actually about 10 minutes was spent in discussing the task to be performed and the purpose of the experiment. Instruction was concluded when the experimenter was satisfied that all important concepts were understood.

A 45-minute practice session preceded the 15 hours of data recording. During this session, the response was continuously monitored and the subject was assured of the quality of his performance. The lack of error feedback made it difficult for the subject to evaluate his own performance until he had some experience with the task.

The instructions were as complete as the subject seemed to need in all but one important area. He was told nothing about the dynamics of the input sequence, except that there would be changes in the probability. He was told to expect both rapid and slow changes. He was instructed that the pay he would receive would be a constant rate per minute of tracking minus the accumulated squared error during the same interval. The amount was computed automatically on an analog computer operating during the tracking session. (The circuit used for the pay scheme is shown in Figure 2.)

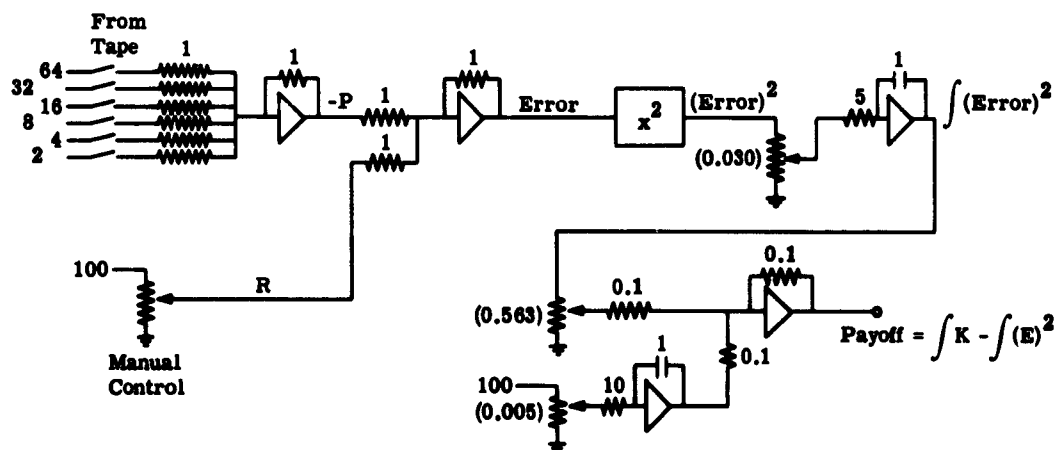


FIGURE 2. ANALOG COMPUTER CIRCUIT FOR PAYOFF

3

EXPERIMENTAL RESULTS

3.1. THE PILOT EXPERIMENT

A pilot experiment was run prior to the main experiment in an attempt to answer three questions. The first concerned the general form and quality of the response. The responses found were qualitatively similar to the response plotted in Figure 3. Both the response to change and the estimation of probability were better than expected.

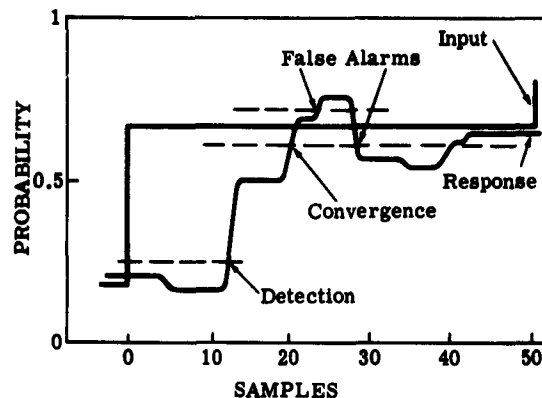


FIGURE 3. A TYPICAL RESPONSE TO A SUBPROBLEM

The second question answered by the pilot experiment concerned changes in response with continued performance of the task, reflecting learning or fatigue. One problem was presented to the four subjects in each of six sessions about two days apart. There was no indication of a significant change in performance after the first session. To test for specific problem learning, the problem which had been presented for six sessions was presented again, but backwards. No decrement in performance was observed. It was concluded that no specific problem learning had occurred. None of the subjects recognized that the problem had been the same in each of the six sessions, nor were they able to describe the changes in the probabilities. Tracking sessions up to 15 minutes caused no particular fatigue or boredom, and it was concluded that sessions of 25 minutes would be permissible on the more isolated, impressive, and comfortable console used in the main experiment.

The third question answered by the pilot experiment concerned the kind and amount of instruction needed to bring the subjects up to a reasonably consistent level of performance. The instructional method described in Section 2 was the result. The subjects in the main experi-

ment performed consistently after the instruction and practice. (Appendix E presents two interesting exceptions.) The important task learning evidently occurs during the first few minutes of performance, and the 45-minute practice session was sufficient.

3.2. RESPONSE MEASURES

The response measures were chosen after study of the Sanborn records from the main experiment. The form of these responses was the same as that in the pilot experiment (shown in Figure 3). The response was characterized by fairly rapid changes separated by periods of little or no change. This discontinuous form indicates that the behavior might be described in terms of a series of decisions concerning changes in the probability. A descriptive model with this characteristic is developed in Section 4. Several of the response measures were chosen to fit this response form. All of the response measures refer to individual subproblems. The following measures were calculated.

(a) DETECTION, D : the number of samples from the step change to the point where the response has changed 0.05 in the direction of the new probability from its value at the point of the step change. If R_n is the response at point n in a subproblem which starts at $n = 1$, the point of detection is the point where $R_n = R_0 \pm 0.05$ (the plus sign indicating an increasing step and the minus a decreasing step).

(b) NO DETECTION, \bar{D} : the number of subproblems in which detection did not occur; that is, R_n never came to within 0.05 of the new probability.

(c) CONVERGENCE, C : the number of samples from the step change to the point where the response is within 0.05 of the new probability. The point of convergence is that at which $R_n = P \pm 0.05$, where P is the probability following the step change. The point of convergence is the first entry into this region from either side.

(d) NO CONVERGENCE, \bar{C} : the number of subproblems in which convergence did not occur; that is, R_n was always outside the 0.05 region about P .

(e) INITIAL CONVERGENCE, IC : the number of subproblems in which the response was within the convergence region about the new probability at the point of the step change. $P - 0.05 \leq R_0 \leq P + 0.05$.

(f) ROOT MEAN SQUARE ERROR, $RMSE$: the square root of the mean square error over the entire subproblem. Error equals the response minus the probability. The response was

measured on a 0 to 1 scale corresponding to the probability measure, and the error can thus be considered an error in probability. For a subproblem of length M,

$$RMSE = \left[\frac{1}{M} \sum_{n=1}^{n=M} (R_n - P)^2 \right]^{1/2}$$

(g) ROOT MEAN SQUARE ERROR AFTER C, $RMSE_C$: the square root of the mean square error from the point of convergence to the end of the subproblem.

$$RMSE_C = \left[\frac{1}{M - C} \sum_{n=C+1}^{n=M} (R_n - P)^2 \right]^{1/2}$$

This measure and the following two were made only when either convergence or initial convergence was measured.

(h) MEAN ERROR AFTER C, ME_C : the mean error from the point of convergence to the end of the subproblem.

$$ME_C = \frac{1}{M - C} \sum_{n=C+1}^{n=M} (R_n - P)$$

(i) FALSE ALARM RATE, FAR: the number of times per sample that the response left the 0.05 convergence region between the point of convergence and the end of the subproblem. If $P - 0.05 \leq R_{n-1} \leq P + 0.05$, and $R_n < P - 0.05$ or $R_n > P + 0.05$, then the point n would be a false alarm point.

Detection and convergence were measures designed to describe the discontinuous response form. The 0.05 criterion used in these measures was selected after an extensive study of the data. In about 80% of the subproblems, a sudden response to the new probability occurred shortly after a step change. This movement was interpreted to be the result of the perception of the change in the probability. The 0.05 detection criterion was selected as measuring this point with fair consistency. For step changes greater than about 0.15, this measure is relatively insensitive to the choice of the 0.05 criterion since the sudden response was characteristically 0.10 or greater.

Convergence is more dependent on the selection of 0.05 as a criterion. The point of convergence was most useful, however, in determining the beginning of measures 7, 8, and 9. These measures were all averaged over flashes, and the location of the convergence point did not affect their values. Detection and convergence, as measured with the 0.05 criterion, are not particularly informative for the smallest step change, 0.06.

Measures 7, 8, and 9, the three starting at the point of convergence, indicate the subject's static estimation ability. The subject is operating under what might be called a dynamic set, however; that is, he has an expectancy for changes in the probability. Changes in his responses during this period could be called microstructure tracking, since the subject was not aware that the probability was constant. No measures were made of the persistence of this microstructure tracking on the longer subproblems. This behavior might begin to diminish with long presentations of a constant probability.

RMSE was the only measure made on all subproblems regardless of their response form. It indicates the overall quality of performance. RMSE is a common measure in continuous tasks of this kind, largely because it is easily derived and manipulated in mathematical expressions.

3.3. DATA ANALYSIS

There was one subproblem for each rate, step size, step direction, probability, constraint, and subject, 3440 in all. The combinations of variables presented here were judged to be the most informative set among the total available from the computer analysis. These quantitative performance measures were the intended output of this experiment, and since no testable hypotheses were generated, no tests of statistical significance were made.

3.4. EXPERIMENTAL DATA

3.4.1. DIFFERENCES BETWEEN SUBJECTS. No qualitative differences existed among the four subjects used in the main experiment. The four subjects in the pilot experiment behaved similarly to those used in the main experiment and to each other. Inspection of the data indicated that for general performance information, it would be best to average the data over subjects. (Appendix F presents some of the subject-by-subject data.)

3.4.2. DETECTION, D. Figures 4 through 7 show the effects of the independent variables on detection. Since the data on step direction show no appreciable difference between positive and negative directions, they are averaged together in all figures. The interaction of step size and rate shown in Figure 4 shows the most interesting relation found. Here detection decreases fairly linearly with step size and increases fairly linearly with the logarithm of rate.

The linear increase in detection with the logarithm of the rate probably reflects a combination of factors influencing the response. A small linear increase with rate would be caused by a constant reaction and movement time. For the usual tracking tasks, this might be expected to be on the order of 0.5 seconds and to yield a lag of 2 flashes at 4 fps and 4 flashes at 8 fps.

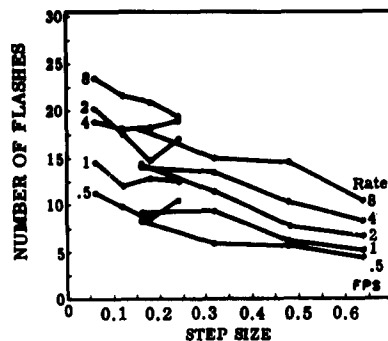


FIGURE 4. DETECTION AS A FUNCTION OF STEP SIZE AND FLASH RATE. Small- and large-step problems are plotted separately, the small extending from 0.06 to 0.24, and the large from 0.16 to 0.64. Detection is measured in flashes.

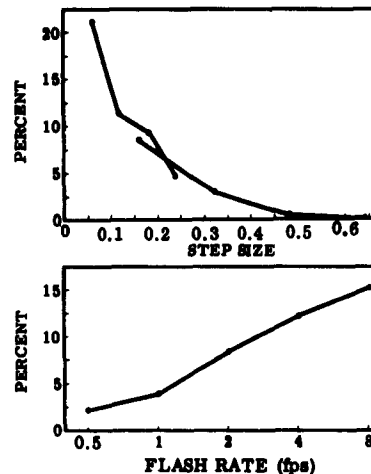


FIGURE 5. PERCENTAGE OF SUBPROBLEMS IN WHICH "NO DETECTION" OCCURS, AS A FUNCTION OF STEP SIZE AND FLASH RATE. The small- and the large-step problems are plotted separately, the small extending from 0.06 to 0.24 and the large from 0.16 to 0.64.

The more important factor is probably a change in the method of performing the task as rate changes. At rates of 0.5 and 1 fps the subjects reported counting the flashes at times, occasionally counting the number of flashes of the lower frequency and comparing this to an estimate of the total number of flashes. They did not use any procedure of this sort consistently, however, at least not one apparent to them. They all reported that the rate of 2 fps was the most difficult. Evidently the methods which they had used effectively at 0.5 and 1 fps became difficult if not impossible at 2 fps. Beginning at 4 fps it is clearly impossible to respond to separate flashes and the series is probably perceived in groups of flashes. The task becomes similar to a continuous tracking task at these rates. Reese [11] postulated that subjects' mechanism for counting light flashes would change at about 4 flashes per second.

Figure 4 shows an effect due to the presentation of the step changes in two separate series, the small- and the large-step problems. There is a region of overlap in step size between these two problems. The smallest change in the large-step problem is 0.16, and the largest in the small-step problem is 0.24. In this overlapping region the small-step problem yields detections of from one to six flashes higher than the large-step problem at all rates. The subjects

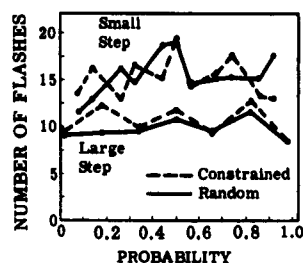


FIGURE 6. DETECTION AS A FUNCTION OF PROBABILITY, CONSTRAINT, AND SMALL- AND LARGE-STEP PROBLEMS. Detection is measured in flashes.

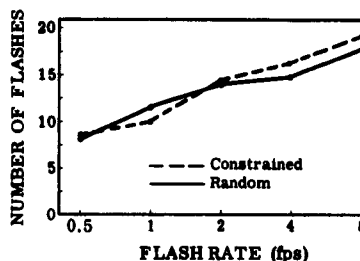
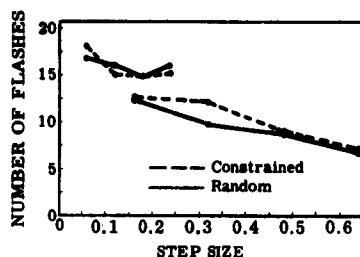


FIGURE 7. DETECTION AS A FUNCTION OF STEP SIZE, SAMPLE RATE, AND CONSTRAINT. The small- and large-step problems are plotted separately, the small extending from 0.06 to 0.24 and the large from 0.16 to 0.64. Detection is measured in flashes.

were evidently modifying their tracking method according to the type of problem being presented. The large- and small-step problems were ordered randomly, of course, and the subjects had no prior indication that there were two problem types. This change is perhaps not surprising considering the difference between the two problem types. The average step changes were 0.15 in the small-step problem and 0.40 in the large-step problem. Step changes of about 0.30 and larger are readily noticed.

The subjects appear to have made larger, more decisive response changes on the large-step problem than on the small-step problem. This more responsive behavior is appropriate in quickly reducing the large errors following the larger step changes.

Figure 5 shows the percentage of "no detections" for the total number of subproblems, as a function of step size and rate. About 90% of the "no detections" occurred with the combination of rate above 4 fps and step size below 0.15. Some of the "no detections" were probably caused by lapses of attention. At 4 fps a 42-flash subproblem is over in 11 seconds. In 25 minutes of continuous tracking a few 11-second lapses are certainly to be expected.

Figure 6 shows the effect of probability on detection. Perhaps the most interesting finding is that detection is not appreciably smaller for the extreme probabilities. In Section 4 it will be seen that responses generated by simple running averages produce detections which are similarly independent of probability.

The variability among a set of detections of a particular step size and rate will depend on the probability, however. Detections of central probabilities, those near 0.5, will have more variability than detections of extreme probabilities, those nearer 0 or 1.

The effect of the flash generation constraint on detection is shown in Figures 6 and 7. Constraint has no particular effect on average detection, but like probability it affects the variability of a set of detections. The constrained problems yield less variable detections.

3.4.3. CONVERGENCE, C. Figures 8 through 12 show the effects of the independent variables on convergence. The interesting effects are again with step size and rate. The effect of rate on convergence is similar to its effect on detection; a linear increase in convergence with the logarithm of rate. The effect of increasing step size is to increase convergence, although the increase is small. The number of flashes between detection and convergence increases as step size increases. This probably reflects the size of the response more than any other factor. Most of the subproblems show a response successively approaching the new probability rather than one that overshoots.

Convergence shows a difference, similar to that noted in detection, between the small- and the large-step problems in the region of overlapping step size.

"No convergence," expressed as a percentage of subproblems, is shown in Figure 9. "No convergence" remains relatively insensitive to changes in step size except for the largest step, 0.64, where it is zero. It is approximately 10% for the large-step problem and 12.5% for the small-step problem. "No convergence" rises sharply with increasing rate, reaching about 28% at 8 fps. This is consistent with the data, which show convergence equal to 35 flashes at 8 fps, about the length of the shortest subproblem.

"Initial convergence" has a high of 35% for a step change of 0.06 and goes to zero for steps of 0.48 and 0.64. It increases slightly with rate from about 8 to 12%.

The relationship between probability and convergence is shown in Figure 11. Convergence is relatively insensitive to probability as was detection.

The effects of constraint on the sample generation are shown in Figures 11 and 12. Again, as with detection, there is little if any effect.

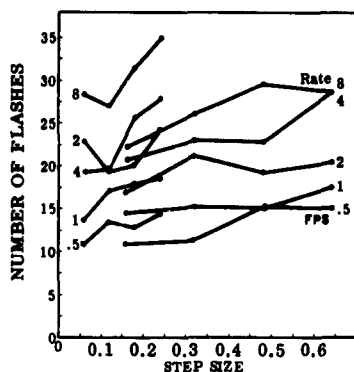


FIGURE 8. CONVERGENCE AS A FUNCTION OF STEP SIZE AND FLASH RATE. The small- and large-step problems are plotted separately, the small extending from 0.06 to 0.24 and the large from 0.16 to 0.64. Convergence is measured in flashes.

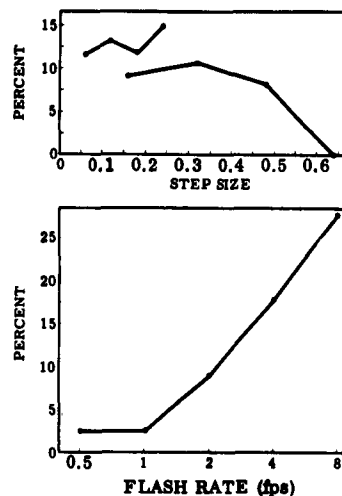


FIGURE 9. PERCENTAGE OF SUBPROBLEMS IN WHICH "NO CONVERGENCE" OCCURS, AS A FUNCTION OF STEP SIZE AND FLASH RATE. The small- and large-step problems are plotted separately, the small extending from 0.06 to 0.24 and the large from 0.16 to 0.64.

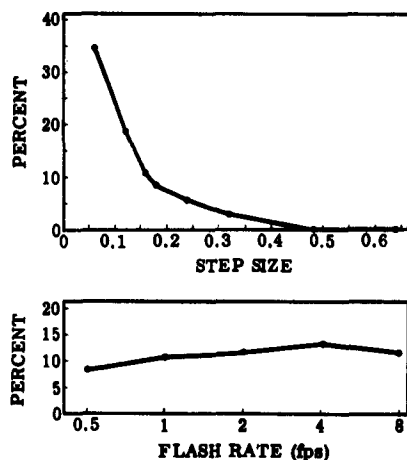


FIGURE 10. PERCENTAGE OF SUBPROBLEMS IN WHICH "ACCIDENTAL INITIAL CONVERGENCE" OCCURS, AS A FUNCTION OF STEP SIZE AND FLASH RATE

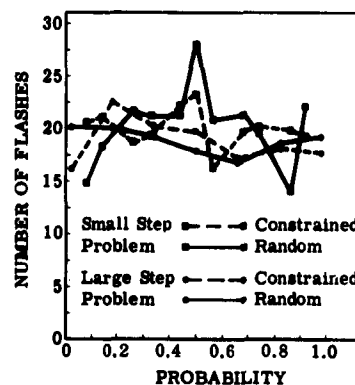


FIGURE 11. CONVERGENCE AS A FUNCTION OF PROBABILITY, CONSTRAINT, AND SMALL- AND LARGE-STEP PROBLEMS. Convergence is measured in flashes.

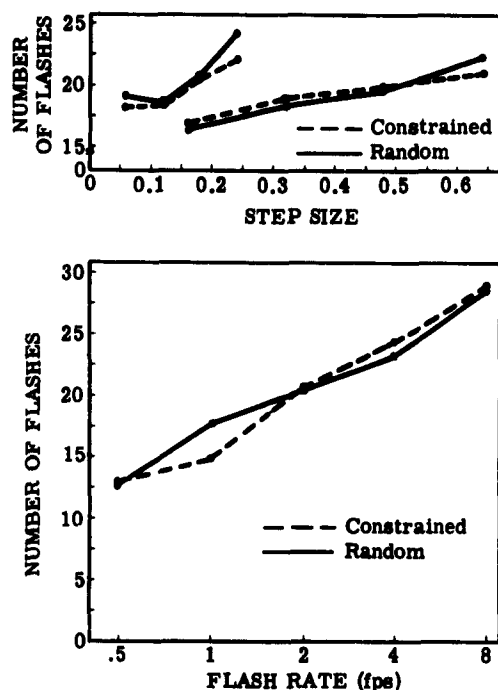


FIGURE 12. CONVERGENCE AS A FUNCTION OF STEP SIZE, FLASH RATE, AND CONSTRAINT. The small- and large-step problems are plotted separately, the small extending from 0.06 to 0.24 and the large from 0.16 to 0.64. Convergence is measured in flashes.

3.4.4. ROOT MEAN SQUARE ERROR, RMSE. This measure was introduced to provide a single overall indicator of the task performance. The most informative variation of RMSE is variation with rate, plotted in Figure 13. RMSE increases linearly with rate from 1 to 8 fps. It is interesting to evaluate this performance measure on a time basis, as might be done when the estimation must take the shortest time possible. Dividing RMSE by fps yields values of error-seconds per flash which decrease as rate increases, going from 0.134 at 1 fps to 0.022 at 8 fps. This decrease might well continue with even higher rates, as the task becomes the tracking of the relative brightness of the lights. Either the limitations on the judgment of relative brightness or simple reaction time would finally limit the performance. This performance index must be viewed with caution. The error itself has a meaningful upper bound at the level

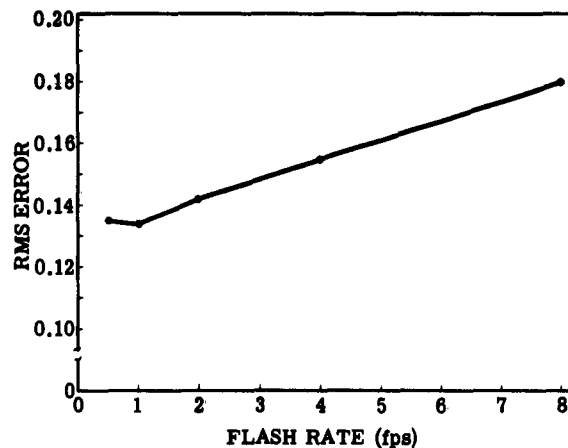


FIGURE 13. ROOT MEAN SQUARE ERROR OVER THE WHOLE SUBPROBLEM AS A FUNCTION OF FLASH RATE

where the lever is left stationary or is moved in some manner independent of the flashes. As this error level is approached, further increases in rate would continue to decrease the index of error-seconds per flash but the index would have little meaning.

The following three measures were made from the point of convergence—or, if initial convergence occurred, from the beginning of the subproblem—to the end of the subproblem. They are therefore measures made on an average of 85 to 90% of all of the subproblems and on about 95% of those subproblems with step changes above 0.15 at rates below 4 fps.

3.4.5. MEAN ERROR AFTER CONVERGENCE, ME_C . The mean error is shown as a function of probability in Figure 14. The average estimate is essentially unbiased at all probabilities. The largest error is smaller than the least scale division on the subject's response indicator, 0.02. Mean error was not significantly affected by rate, constraint, step size, or subjects.

This finding contradicts a body of conjecture based in part on the results of static estimation and choice experiments. Neither the overestimation of high nor the underestimation of low probabilities appears. The excellence in static estimation was undoubtedly due at least in part to the two distinctive features of the task, the dynamic estimation and the use of the tracking lever as the response mechanism.

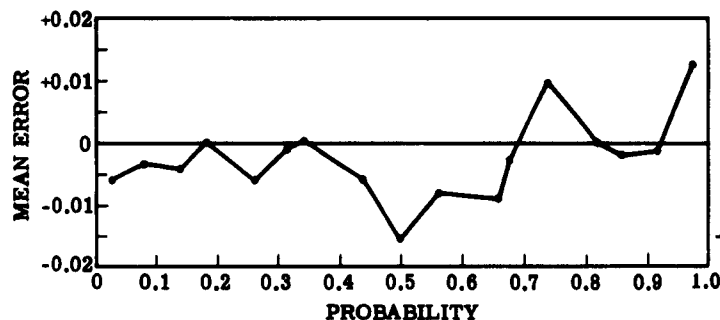


FIGURE 14. MEAN ERROR AS A FUNCTION OF PROBABILITY. This measure is made from convergence to the end of the subproblem.

3.4.6. ROOT MEAN SQUARE ERROR AFTER CONVERGENCE, $RMSE_C$. $RMSE_C$ is shown in Figures 15 through 17. The only independent variable not affecting $RMSE_C$ is step size. This indicates that the period after the point of convergence is not affected by step size. The constraint on the generation of the flashes reduced the $RMSE_C$ by about 0.014 and does not appear to interact with either step size or rate. $RMSE_C$ decreases with increasing rate from 0.5 to 2 fps and thereafter remains relatively constant.³ Considered together with the data indicating smaller detection values at the lower rates, it is highly probable that the number of decisions concerning changes in the probability on a per flash basis is highest at the lowest rate. Thus the additional decision time available at the lower rates permitted smaller detection values but resulted in larger $RMSE_C$ when the probability was constant.

The effect of probability on the $RMSE_C$ is shown in Figure 17. The "random" problems are consistently higher than the "constrained" problems at all probabilities. The $N = 17.3$ line is the $RMSE$, or standard deviation, of a 17.3 flash average. The subject's response is about this good or better at all probabilities.

3.4.7. FALSE ALARM RATE, FAR. The number of false alarms per flash is shown in Figures 18 through 20. Its behavior is similar to $RMSE_C$. It is similarly insensitive to the size of the step change. Increasing rate causes a decrease in FAR up to 4 fps with an apparent leveling off above 4 fps. These data lend additional support to the hypothesis concerning an increase in number of decisions per flash at the lower rates. False alarms can be considered as indicating decisive changes in the estimate.

³Variation among subjects was high for 4 and 8 fps. See Figure 28, Appendix F.

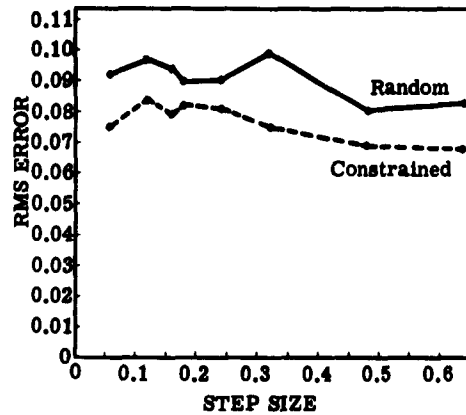


FIGURE 15. ROOT MEAN SQUARE ERROR AS A FUNCTION OF STEP SIZE AND CONSTRAINT. This measure is made from convergence to the end of the subproblem.

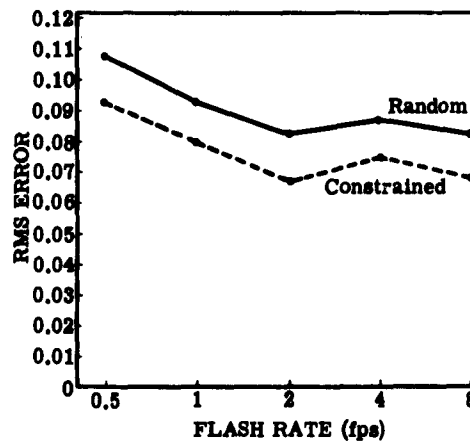


FIGURE 16. ROOT MEAN SQUARE ERROR AS A FUNCTION OF FLASH RATE AND CONSTRAINT. This measure is made from convergence to the end of the subproblem.

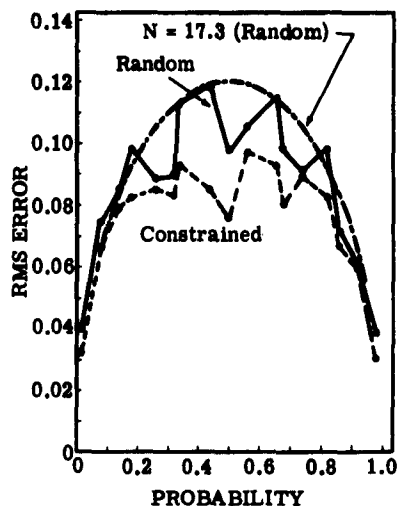


FIGURE 17. ROOT MEAN SQUARE ERROR AS A FUNCTION OF PROBABILITY AND CONSTRAINT. This measure is made from convergence to the end of the subproblem. The standard deviation for a 17.3 sample mean is shown for the random problem.

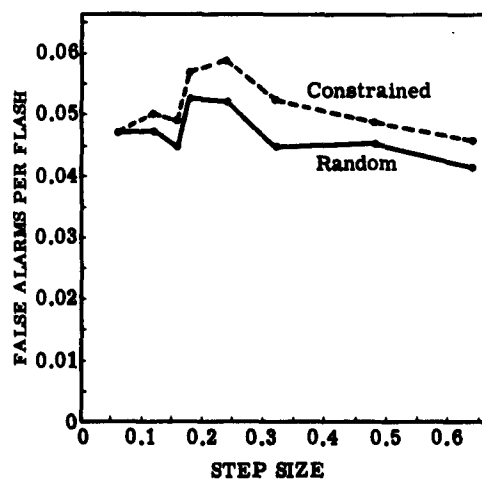


FIGURE 18. FALSE ALARM RATE IN FALSE ALARMS PER FLASH AS A FUNCTION OF STEP SIZE AND CONSTRAINT. This measure is made from convergence to the end of the subproblem.

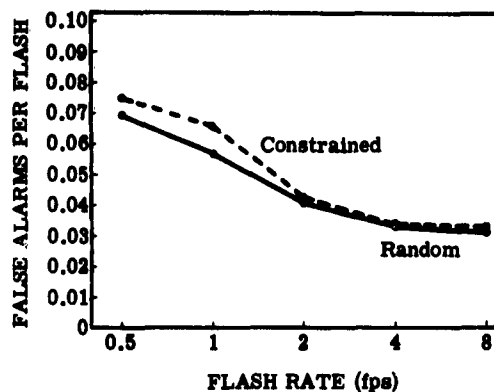


FIGURE 19. FALSE ALARM RATE IN FALSE ALARMS PER FLASH AS A FUNCTION OF FLASH RATE AND CONSTRAINT. This measure is made from convergence to the end of the subproblem.

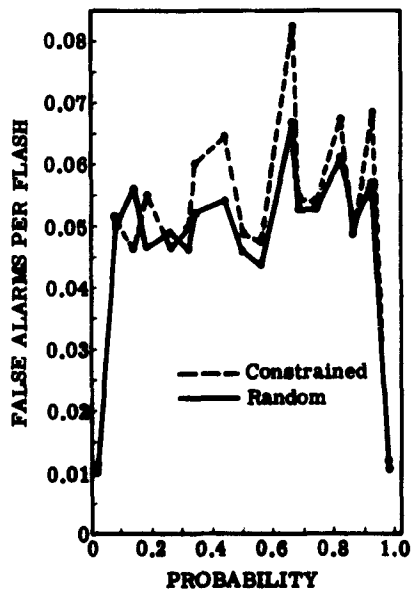


FIGURE 20. FALSE ALARM RATE IN FALSE ALARMS PER FLASH AS A FUNCTION OF PROBABILITY AND CONSTRAINT. The measure is made from convergence to the end of the subproblem.

FAR remains constant over the entire probability range with the exception of the extreme values, 0.02 and 0.98. These probabilities were usually estimated as 0 or 1, with an excursion away from 0 or 1 only after an occurrence of the infrequent flash. Since FAR did not change with probability, it appears that the rate of "decisive" movements (greater than 0.05 from the probability) remained constant for all probabilities. The reduction in $RMSE_C$ as the probability tends to extreme values therefore indicates that the time spent at these "erroneous" estimates decreased with extreme probabilities. This hypothesis is supported by observations made during the tracking sessions. The lever movements appeared larger although less frequent at the more extreme probabilities. The increase in magnitude evidently compensated for the decrease in frequency to maintain the FAR at a constant level.

The constrained series produced a slightly higher false alarm rate than the random series. The constrained series has a greater number of runs of right or left flashes and would be expected to yield a higher decision rate. All of the FAR data will be dependent on the false alarm criterion level. A larger criterion could well reverse the constraint finding, for example, since the random series probably produces larger decision movements than the constrained series.

3.5. SUMMARY OF RESULTS

The response to a step change in probability can be described in three regions: the period before any response to the change, before the point of detection; the period before the convergence on a new estimate; and the period from the convergence point to the end of the subproblem. These regions were defined mathematically as functions of probability response form and somewhat arbitrary constants in order to achieve a complete description of the response.

Detection increases with increasing rate and decreases with increasing step size. The range was from 4 to 24 flashes for a rate range of 0.5 to 8 flashes per second and a step size range of 0.06 to 0.64. Detection was approximately nine flashes for a step of 0.32 at 1 fps.

Convergence increases with both rate and step size. The range was from 11 to 35 flashes for the same step and rate ranges stated above. Convergence was approximately 15 flashes for a step of 0.32 at 1 fps.

Both detection and convergence were independent of the constraint imposed on the generation of the flash series. Both were independent of probability.

After the point of convergence the average estimate was unbiased at all probabilities. This unbiased estimate had an RMS error, or standard deviation, of about 0.06.

The overall task performance was measured by the RMS error throughout the subproblem. RMSE increased linearly with rate from 0.135 at 1 fps to 0.180 at 8 fps.

4

MATHEMATICAL MODELS

Three mathematical models will be derived in this chapter. Two of these will be called normative, since the purpose for their derivation is to provide standards with which to compare the data presented in Section 3. The third is a descriptive model designed to simulate the human performance.

The somewhat arbitrary forms of the normative models, and the optimized parameters used, were selected to provide the best RMS error fit to the various inputs used in the probability tracking task. One form selected is a constant weighted average over a finite number of past flashes. The simplicity of this model makes it ideal for intuitive comparisons with the subjects' performance. The number of flashes in the running average is selected to give the best fit. The other model has geometrically decreasing weight for flashes extending into the past. This model is more appealing from the standpoint of response to the step inputs. It also corresponds to assumptions often made concerning the human immediate memory function. The best fit is found by selecting the appropriate geometric ratio.

More sophisticated linear models, and certainly some nonlinear ones, would undoubtedly perform this task with a lower RMSE than the two models selected. The value of more complex models for providing simple standards is marginal, however.

The descriptive model was derived from thoughts on how the subjects performed the task. Its form arises from the qualitative aspects of the data and from observations of the subjects' behavior. It has four parameters, which are adjusted to yield a minimum RMSE fit to a subject's response.

The normative models to be considered have the form

$$r(n) = \sum_{i=1}^{i=N} w_i s_{n-i+1} \quad (1)$$

where $r(n)$ is the model's response or output at the point n in the sample series, and w_i is a weight attached to the sample s_{n-i+1} . The response at n is thus the weighted average of the sample at n and its $N - 1$ immediate predecessors. This is an averaging or smoothing model intuitively appropriate to this task. It is limited to samples at and prior to the response point, considering only a finite number of these, and is therefore physically realizable. w_i is not a function of n and could be described as sample-invariant.

The random variables s_n are drawn from an infinite population and are independent. They have values 0 or 1, corresponding respectively to left and right on the subject's display. The probability of a 1 is P , and the probability of a 0 is therefore $1 - P$.

When the N samples are all generated from a static distribution described by the probability P , it is desirable that the estimate be unbiased or that

$$\overline{r(n)} = P \quad (2)$$

where $\overline{r(n)}$ is the expected value of $r(n)$, an ensemble average. This simply requires that

$$\sum_{i=1}^{i=N} w_i = 1 \quad (3)$$

The responses of the two model forms will be derived for a subproblem beginning with $n = 1$ as the first sample of the new probability and ending with $n = M$. The previous probability will be P_1 and the subproblem probability P_2 . The step change is therefore $P_2 - P_1$. For $N < n \leq M$ the response will be called steady-state, since the samples are all from a static distribution. For $1 \leq n \leq N$ the response will be called transient.

4.1. A MODEL WITH GEOMETRIC WEIGHTING

The first model to be considered has a weighting function

$$w_i = ar^{i-1} \quad (4)$$

where a and r are constants and $0 < r < 1$. This function assigns geometrically decreasing weights to the samples. Limiting r to the range 0 to 1 confines the function to one assigning monotonically decreasing weights to samples receding from n .

The value of N , the number of samples included in one computation, will be selected as a number large enough to assure the relative unimportance of the weight at $n = N$, ar^{N-1} , compared to the weight at $n = 1$, a . This merely implies that the function's memory extends smoothly to the point of essentially complete "forgetting." The exact value of N in any particular model of this form is relatively unimportant to the considerations that follow. It will simply be assumed that

$$r^N \ll 1 \quad (5)$$

and all quantities of this magnitude will be dropped.

Of primary interest is the selection of r to produce an optimum model, that is, one having the least mean squared error. This particular measure of performance was the same one used in measuring of the subject's performance and in the payoff scheme.

The constant a is selected to satisfy Equation 3, which becomes

$$\sum_{i=1}^{i=N} ar^{i-1} = a + ar + ar^2 + \dots + ar^{N-1} = a \frac{1-r^N}{1-r} = 1 \quad (6)$$

or

$$a = 1 - r, r^N \ll 1 \quad (7)$$

We will be concerned with two quantities: $\overline{r(n)}$, the expected value of $r(n)$ at the point n , and $\sigma_r^2(n)$, the variance of $r(n)$ at the point n . These are ensemble averages.

For the expected value of r we have

$$\overline{r(n)} = E \left[\sum_{i=1}^{i=N} w_i s_{n-i+1} \right] \quad (8)$$

and since the w_i are constant over the ensemble,

$$\overline{r(n)} = \sum_{i=1}^{i=N} w_i E(s_{n-i+1}) = \sum_{i=1}^{i=N} w_i \quad (9)$$

For the step function input, $\overline{r(n)}$ will depend on P_1 and P_2 during the transient phase and on P_2 along during the steady-state phase.

For $1 \leq n \leq N$ we have

$$\begin{aligned} r(n) &= \frac{(1-r)(1-r^n)}{1-r} P_2 + \left[1 - \frac{(1-r)(1-r^n)}{1-r} \right] P_1 \\ &= (1-r^n) P_2 + r^n P_1 \\ &= P_2 - r^n (P_2 - P_1) \end{aligned} \quad (10)$$

For $N < n \leq M$,

$$\overline{r(n)} = P_2 \quad (11)$$

For the variance we have the variance of the sum of the $w_i s_{n-i+1}$ terms. Since the s_n are independent, we have the sum of the variances of the individual terms

$$\sigma_r^2(n) = a^2 \sigma_s^2(n) + a^2 r^2 \sigma_s^2(n-1) + a^2 r^4 \sigma_s^2(n-2) + \dots + a^2 r^{2(N-1)} \sigma_s^2(n-N+1) \quad (12)$$

where $\sigma_s^2(n)$ is the variance of the sample s_n . Since $\sigma_s^2(n) = P(1 - P)$, where P is the probability with which s_n was generated, $\sigma_s^2(n)$ will be a constant for constant P , and in particular it will have two values during the transient phase, $\sigma_1^2 = P_1(1 - P_1)$ and $\sigma_2^2 = P_2(1 - P_2)$.

For the transient phase we then have

$$\begin{aligned}\sigma_r^2(n) &= (1-r)^2 \frac{(1-r^{2n})}{(1-r^2)} \sigma_2^2 + \left[(1-r)^2 \frac{(1-r^{2N})}{(1-r^2)} - (1-r)^2 \frac{(1-r^{2n})}{(1-r^2)} \right] \sigma_1^2 \\ &= \frac{(1-r)}{(1+r)} \left[\sigma_2^2 - r^{2n} (\sigma_2^2 - \sigma_1^2) \right], \quad r^N < 1\end{aligned}\quad (13)$$

and for the steady-state phase,

$$\sigma_r^2(n) = \frac{(1-r)}{(1+r)} \sigma_2^2 \quad (14)$$

We can now proceed with the formulation of the model's performance in terms of its mean square error. The error during the transient phase can be written as

$$\begin{aligned}e(n) &= r(n) - P_2 \\ &= [\bar{r}(n) - P_2] + [\bar{r}(n) - r(n)]\end{aligned}\quad (15)$$

and the squared error is then

$$e^2(n) = [\bar{r}(n) - P_2]^2 + [r(n) - \bar{r}(n)]^2 + 2[\bar{r}(n) - P_2][r(n) - \bar{r}(n)] \quad (16)$$

The expected value of the squared error is then

$$\overline{e^2(n)} = [\bar{r}(n) - P_2]^2 + \sigma_r^2(n) \quad (17)$$

since $[\bar{r}(n) - P_2]$ is a constant for a particular n and $E[r(n) - \bar{r}(n)]$ is 0.

The average value of this mean square error over the transient phase of the subproblem will be $\langle e^2 \rangle_T$, representing the average over the ensemble and also over samples in the subproblem. We will then have

$$\langle e^2 \rangle_T = \frac{1}{N} \sum_{n=1}^{n=N} \overline{e^2(n)} = \frac{1}{N} \sum_{n=1}^{n=N} [\bar{r}(n) - P_2]^2 + \frac{1}{N} \sum_{n=1}^{n=N} \sigma_r^2(n) \quad (18)$$

When Equation 10 is used, the first term on the right side of Equation 18 becomes

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^{n=N} [\bar{r}(n) - P_2]^2 &= \frac{1}{N} \sum_{n=1}^{n=N} [r^n(P_1 - P_2)]^2 \\ &= \frac{(P_1 - P_2)^2}{N} \frac{r^2}{1 - r^2}, r^N \ll 1 \end{aligned} \quad (19)$$

When Equation 13 is used, the second term on the right side of Equation 18 becomes

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^{n=N} \sigma_r^2(n) &= \frac{1}{N} \sum_{n=1}^{n=N} \frac{(1-r)}{(1+r)} \left[\sigma_2^2 - r^{2n} (\sigma_2^2 - \sigma_1^2) \right] \\ &= \frac{(1-r)}{N(1+r)} \left[N\sigma_2^2 - \frac{r^2}{1-r^2} (\sigma_2^2 - \sigma_1^2) \right] \\ &= \frac{(1-r)}{(1+r)} \left[\sigma_2^2 - \frac{r^2}{N(1-r^2)} (\sigma_2^2 - \sigma_1^2) \right], r^N \ll 1 \end{aligned} \quad (20)$$

The average mean square error during the transient phase, $\langle \bar{e}^2 \rangle_T$, is then

$$\langle \bar{e}^2 \rangle_T = \frac{(P_1 - P_2)^2}{N} \frac{r^2}{1 - r^2} + \frac{(1-r)}{(1+r)} \left[\sigma_2^2 - \frac{r^2}{N(1-r^2)} (\sigma_2^2 - \sigma_1^2) \right] \quad (21)$$

The average mean square error during the steady-state phase $\langle \bar{e}^2 \rangle_{SS}$ is simply the variance $\sigma_r^2(n)$, is given by Equation 14:

$$\langle \bar{e}^2 \rangle_{SS} = \frac{(1-r)}{(1+r)} \sigma_2^2 \quad (22)$$

The average mean square error over the whole subproblem is then

$$\begin{aligned} \langle \bar{e}^2 \rangle_{SP} &= \frac{N}{M} \langle \bar{e}^2 \rangle_T + \frac{M-N}{M} \langle \bar{e}^2 \rangle_{SS} \\ &= \frac{(P_1 - P_2)^2}{M} \frac{r^2}{1 - r^2} + \frac{N}{M} \frac{(1-r)}{(1+r)} \left[\sigma_2^2 - \frac{r^2}{N(1-r^2)} (\sigma_2^2 - \sigma_1^2) \right] \\ &\quad + \left(\frac{M-N}{M} \right) \frac{(1-r)}{(1+r)} \sigma_2^2 \\ &= \frac{(P_1 - P_2)^2}{M} \frac{r^2}{(1-r^2)} + \frac{(1-r)}{(1+r)} \left[\sigma_2^2 - \frac{r^2}{M(1-r^2)} (\sigma_2^2 - \sigma_1^2) \right] \end{aligned} \quad (23)$$

We are interested in the performance of this model over the types of problem given to the subjects. The average mean square error over a problem is given by

$$\overline{\langle e^2 \rangle}_P = \sum_{i=1}^{i=S} \frac{M_i}{T} \left\{ \frac{(P_{i-1} - P_i)^2}{M_i} \frac{r^2}{1-r^2} + \frac{1-r}{1+r} \left[\sigma_i^2 - \frac{r^2}{M_i(1-r^2)} (\sigma_i^2 - \sigma_{i-1}^2) \right] \right\} \quad (24)$$

where M_i is the length of subproblem i , T is the total problem length in samples, and S is the total number of subproblems.

This expression can be simplified by making the following assumptions based on the methods used for generating the problems (Appendix A). The M_i were selected randomly, without replacement, from a set of equally frequent values and assigned to the subproblems. The sum of $(\sigma_i^2 - \sigma_{i-1}^2)/M_i$ therefore approaches zero for long series of subproblems. Similarly the term $M_i \sigma_i^2/T$ will approach simply σ_i^2/S . Equation 24 can therefore be written as

$$\overline{\langle e^2 \rangle}_P = \frac{r^2}{1-r^2} \frac{1}{T} \sum_{i=1}^{i=S} (P_{i-1} - P_i)^2 + \frac{1-r}{1+r} \frac{1}{S} \sum_{i=1}^{i=S} \sigma_i^2 \quad (25)$$

The large-step problem had values of $|P_1 - P_2|$ of 0.16, 0.32, 0.48, and 0.64, occurring in 12, 10, 8, and 6 subproblems respectively, yielding

$$\frac{1}{T} \sum_{i=1}^{i=S} (P_{i-1} - P_i)^2 = 0.00251$$

σ_i^2 had values of 0.250, 0.224, 0.148, and 0.020 occurring in 6, 12, 10, and 8 subproblems respectively, yielding

$$\frac{1}{S} \sum_{i=1}^{i=S} \sigma_i^2 = 0.162$$

T was 2241 samples and S was 36 subproblems. The average mean square error over the large-step problem is then

$$\overline{\langle e^2 \rangle}_{LSP} = 0.00251 \frac{r^2}{1-r^2} + 0.162 \frac{1-r}{1+r} \quad (26)$$

The small-step problem had values of $|P_1 - P_2|$ of 0.06, 0.12, 0.18, and 0.24, occurring in 12, 10, 16, and 12 subproblems respectively, and yielding

$$\frac{1}{T} \sum_{i=1}^{i=S} (P_{i-1} - P_i)^2 = 0.000465$$

σ_1^2 had values of 0.250, 0.246, 0.217, 0.192, 0.120, and 0.085, occurring in 6, 10, 12, 10, 6, and 6 subproblems respectively, and yielding

$$\frac{1}{S} \sum_{i=1}^{i=S} \sigma_i^2 = 0.194$$

T was 3000 samples and S was 50 subproblems. The average mean square error over the small-step problem is then

$$\overline{\langle e^2 \rangle}_{SSP} = 0.000465 \frac{r^2}{1-r^2} + 0.194 \frac{1-r}{1+r} \quad (27)$$

It will be of interest, for comparative purposes, to evaluate this model for the case in which only one value of r is used for both the large- and small-step problems. This model will be called nondiscriminating in Section 5. In this case the sums in Equation 25 are over both problem types, with T being equal to 5241 samples and S being 86 subproblems. We have

$$\frac{1}{T} \sum_{i=1}^{i=S} (P_{i-1} - P_i)^2 = 0.00134$$

and

$$\frac{1}{S} \sum_{i=1}^{i=S} \sigma_i^2 = 0.181$$

The average mean square error over the large- plus the small-step problems is then

$$\overline{\langle e^2 \rangle}_{S+L} = 0.00134 \frac{r^2}{1-r^2} + 0.181 \frac{1-r}{1+r} \quad (28)$$

We are interested in the selection of an optimum value of r for these three problem types. Equation 25 can be written as

$$\overline{\langle e^2 \rangle}_P = k \frac{r^2}{1-r^2} + v \frac{1-r}{1+r} \quad (29)$$

where k and v are the constants for the specific problem type. The minimum of this function over r can then be found by setting

$$\frac{d \overline{\langle e^2 \rangle}_P}{dr} = \frac{-2vr^2 + (2k + 4v)r - 2v}{(1-r^2)^2} = 0 \quad (30)$$

which yields two roots

$$r_{1,2} = \frac{k+2v}{2v} \pm \left[\left(\frac{k+2v}{2v} \right)^2 - 1 \right]^{1/2} \quad (31)$$

The minus sign yields a value of r between 0 and 1 and is also the minimum.

For the large-step problem Equation 31 gives an optimum $r = 0.883$. Using this value of r in Equation 26 we have a corresponding minimum mean square error of 0.0190.

For the small-step problem Equation 31 gives an optimum $r = 0.953$. Using this value of r in Equation 27, we have a corresponding minimum mean square error of 0.00931.

For the large- plus small-step problems Equation 31 gives an optimum $r = 0.915$. Using this value of r in Equation 28, we have a corresponding minimum mean square error of 0.0147.

4.2. A MODEL WITH CONSTANT WEIGHTING

The second model gives a constant weight to each of N samples; that is, it is a simple averaging model. The derivation of the response and errors for this model will parallel that for the geometric model, and some of the detailed explanations will be omitted.

The weighting function is

$$w_i = 1/N \quad (32)$$

where N is the number of samples in the average and the weight is $1/N$ to satisfy Equation 3.

In this case the transient response will be

$$\overline{r(n)} = \frac{n}{N} P_2 + \frac{N-n}{N} P_1 = P_1 + \frac{P_2 - P_1}{N} n \quad (33)$$

and the steady-state response

$$\overline{r(n)} = P_2 \quad (34)$$

The variance of $r(n)$ during the transient phase will be the variance of the sum of N terms, each with weight $1/N$:

$$\begin{aligned} \sigma_r^2(n) &= \frac{N-n}{N^2} \sigma_1^2 + \frac{n}{N^2} \sigma_2^2 \\ &= \frac{1}{N} \left[\sigma_1^2 + \frac{n}{N} (\sigma_2^2 - \sigma_1^2) \right] \end{aligned} \quad (35)$$

The variance of $r(n)$ during the steady-state phase will simply be

$$\sigma_r^2(n) = \frac{\sigma_2^2}{N} \quad (36)$$

Following the same procedures and arguments developed in the derivation of Equations 15 through 19, we have

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^{n=N} [\overline{r(n)} - P_2]^2 &= \frac{1}{N} \sum_{n=1}^{n=N} \left[P_1 + \frac{P_2 - P_1}{N} n - P_2 \right]^2 \\ &= \frac{(P_2 - P_1)^2}{N} \sum_{n=1}^{n=N} \left(\frac{n}{N} - 1 \right)^2 \\ &= (P_2 - P_1)^2 \left(\frac{1}{3} - \frac{1}{2N} + \frac{1}{6N^2} \right) \end{aligned} \quad (37)$$

and analogous to Equation 20 we have

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^{n=N} \sigma_r^2(n) &= \frac{1}{N^2} \sum_{n=1}^{n=N} \left[\sigma_1^2 + \frac{n}{N} (\sigma_2^2 - \sigma_1^2) \right] \\ &= \frac{1}{N} \sigma_1^2 + \frac{N+1}{2N^2} (\sigma_2^2 - \sigma_1^2) \end{aligned} \quad (38)$$

The average mean square error during the transient phase is then

$$\overline{e^2}_T = (P_2 - P_1)^2 \left(\frac{1}{3} - \frac{1}{2N} + \frac{1}{6N^2} \right) + \frac{1}{N} \sigma_1^2 + \frac{N+1}{2N^2} (\sigma_2^2 - \sigma_1^2) \quad (39)$$

The average mean square error during the steady-state phase is the variance, given by Equation 36.

$$\overline{e^2}_{SS} = \frac{\sigma_2^2}{N} \quad (40)$$

The average mean square error over the whole subproblem is then

$$\begin{aligned} \overline{e^2}_{SP} &= \frac{N}{M} \overline{e^2}_T + \frac{M-N}{M} \overline{e^2}_{SS} \\ &= \frac{N}{M} (P_1 - P_2)^2 \left(\frac{1}{3} - \frac{1}{2N} + \frac{1}{6N^2} \right) + \frac{1-N}{2MN} (\sigma_2^2 - \sigma_1^2) + \frac{\sigma_2^2}{N} \end{aligned} \quad (41)$$

As with the geometric weighted model, we are interested in the performance of this model over the problems given to the subjects. The average mean square error over a problem is given by

$$\langle e^2 \rangle_P = \sum_{i=1}^{i=S} \frac{M_i}{T} \left\{ \frac{(P_{i-1} - P_i)^2}{M_i} \left(\frac{N}{3} - \frac{1}{2} + \frac{1}{6N} \right) + \frac{1-N}{2M_i N} (\sigma_i^2 - \sigma_{i-1}^2) + \frac{\sigma_i^2}{N} \right\} \quad (42)$$

Using the same arguments as those leading to Equation 25, we have

$$\langle e^2 \rangle_P = \frac{1}{T} \left(\frac{N}{3} - \frac{1}{2} + \frac{1}{6N} \right) \sum_{i=1}^{i=S} (P_{i-1} - P_i)^2 + \frac{1}{NS} \sum_{i=1}^{i=S} \sigma_i^2 \quad (43)$$

These sums are the same as those calculated for the geometric weighted model. For the large-step problem we have

$$\langle e^2 \rangle_{LSP} = 0.00251 \left(\frac{N}{3} - \frac{1}{2} + \frac{1}{6N} \right) + 0.162 \frac{1}{N} \quad (44)$$

For the small-step problem,

$$\langle e^2 \rangle_{SSP} = 0.000465 \left(\frac{N}{3} - \frac{1}{2} + \frac{1}{6N} \right) + 0.194 \frac{1}{N} \quad (45)$$

And for the large- plus small-step problems,

$$\langle e^2 \rangle_{L+S} = 0.00134 \left(\frac{N}{3} - \frac{1}{2} + \frac{1}{6N} \right) + 0.181 \frac{1}{N} \quad (46)$$

The minima can again be selected by letting k and v be the constants for the particular problem type

$$\langle e^2 \rangle_P = k \left(\frac{N}{3} - \frac{1}{2} + \frac{1}{6N} \right) + \frac{v}{N} \quad (47)$$

and solving

$$\frac{d \langle e^2 \rangle_P}{dN} = \frac{k}{3} - \frac{k}{6N^2} - \frac{v}{N^2} = 0 \quad (48)$$

This yields

$$N = \left[\frac{1}{2} + \frac{3v}{k} \right]^{1/2} \quad (49)$$

For the large-step problem Equation 49 gives an optimum $N = 14.1$. Using this value of N in Equation 44 we have a corresponding minimum mean square error of 0.0220.

For the small-step problem Equation 49 gives an optimum $N = 35.3$. Using this value of N in Equation 45, we have a corresponding minimum mean square error of 0.0107.

For the large- plus small-step problems, Equation 49 gives an optimum $N = 20.1$. Using this value of N in Equation 46 we have a corresponding minimum mean square error or 0.0172.

4.3. A DESCRIPTIVE MODEL

Inspection of data on the subjects' response shows that they did not perform the estimation task as smoothly as the two normative models. The responses were characterized by rapid adjustments separated by periods of little or no movement. This evidence, together with thoughts on how this task might be performed, led to the postulation of the following model as an attempt to describe the human performance.

This model operates as follows. The subject maintains a short running average of the previous k_3 flashes. This average is of exactly the same type as the second normative model discussed above. At each flash this average is compared with the existing setting of the response lever and the difference noted. If this distance measure is greater than a prescribed criterion level, the response is changed to a new value at some point intermediate between the old response and the running average. If the difference is less than the criterion level, the response remains unchanged.

Several features make this descriptive model attractive. It uses the lever as a memory device, moving it only a fraction of the distance to the new average and thus preserving some of the information in the previous setting. This memory function permits a smaller number of flashes in the running average than would otherwise be required to produce the levels of mean square error measured from the subject's responses. The criterion level corresponds to the concept of the subjects' smallest perceptible difference between the running average and the lever position. It permits the response to remain stationary during periods when the running average deviates only slightly from the response.

This model's operation can be thought of as a form of hypothesis testing. At each flash it is testing the hypothesis that the running average is from a population described by the response lever setting, using the criterion level as a form of significance measure. The subject's performance is thus viewed as a succession of decision making situations. This framework is appropriate to the inclusion of more higher mental processes than are in the usual manual tracking task.

This model can be described mathematically as follows:

$$u(n) = \frac{1}{k_3} \sum_{i=1}^{i=k_3} s_{n-i+1} \quad (50)$$

where $u(n)$ is the running average of k_3 flashes, s_n . If $r(n)$ is the current lever setting and

$$|r(n) - u(n)| \leq k_1 \quad (51)$$

where k_1 is the criterion band, then

$$r(n+1) = r(n) \quad (52)$$

If, however,

$$|r(n) - u(n)| > k_1 \quad (53)$$

then

$$r(n+1) = r(n) + k_2 [u(n) - r(n)] \quad (54)$$

where k_2 is the fractional lever adjustment.

A fourth parameter, k_4 , was also considered. It represents a time (flash) shift between the subject's and the model's responses. The subject's response at n was compared to the model's at $n - k_4$.

The four parameters are constrained to the following ranges:

$$0 < k_1 < 1 \quad (55)$$

where 0 yields adjustment decisions at each sample and 1 yields no adjustment decisions.

$$0 < k_2 \leq 1 \quad (56)$$

where 0 yields no response changes and 1 represents simply the following of the running average whenever an adjustment decision is made:

$$1 \leq k_3 \leq K, \quad k_3 \text{ an integer} \quad (57)$$

where K is some reasonable maximum number of flashes that the subject could be expected to assimilate in one averaging calculation. No definite values for K are known for this task. It is certainly reasonable to assume that the flashes are not simply remembered as a succession of binary symbols but are encoded into a larger symbol set; perhaps one depending on the lengths of runs of one of the binary symbols. Considering the nature of the task and its difficulty, it

would seem unlikely that more than 20 flashes could be used in an averaging calculation; a value closer to 10 would be more appropriate.

There are no particular constraints on k_4 except that $k_4 < 0$ implies subject prediction with respect to the model.

Once the form of the model has been thus chosen, the task is to select parameter sets (k_1, k_2, k_3, k_4) which will make the model best describe the human performance. The criterion used for this selection was the minimization of the mean square error between the subject's and the model's responses over particular problems. This measure was selected as providing the best signal measure of performance, as was mentioned in Section 2. The selection of the minimum mean square error for the criterion assures a fairly close fit to the transient portion of the response where the error is large, at the possible expense of fit to the steady-state portion.

The actual minimization process was carried out as follows. The model was programmed on an IBM 709 computer. The computer was then fed four of the input problems used in the experiment plus the responses of one of the four subjects to these problems. At each sample point the squared difference between the subject's and the model's responses was calculated and accumulated. The values obtained were simply printed out at the end of each problem-parameter set combination. The large number of parameter sets possible and the possibility of numerous minima precluded the use of an automatic searching technique for the minima. Several computer runs were made in which previously selected parameter ranges were either extended or filled in, according to the results of the previous run. The total variation of the parameters was through the following ranges.

$$0.02 \leq k_1 \leq 0.20 \quad (6 \text{ values})$$

$$0.10 \leq k_2 \leq 0.90 \quad (7 \text{ values})$$

$$1 \leq k_3 \leq 28 \quad (12 \text{ values})$$

$$-2 \leq k_4 \leq 4 \quad (6 \text{ values})$$

The four problems investigated were the large- and small-step problems, random constraint, at 1 and 4 fps. The subject was S-2.

Several parameter sets with approximately equal minimum error measures were found for each problem type. In each case these minima represented either a valley in the error function or fairly distinct minima separated by regions of higher error. Table I shows the various parameter sets and their corresponding minimum errors, $\langle e_{MS}^2 \rangle$. In each group of parameter sets one can find various tradeoffs among the parameters which yield the approximately equal error measures.

TABLE I. PARAMETER SETS FOR THE DESCRIPTIVE MODEL YIELDING
MINIMUM VALUES OF $\langle e_{MS}^2 \rangle_P$

	Criterion k_1	Fractional Adjustment k_2	Memory k_3	Lag k_4	$\langle e_{MS}^2 \rangle_P$	$\langle e_M^2 \rangle_P$	$\langle 2e_M e_{MS} \rangle_P$	ρ
Large Step, 1 fps								
1.	0.05	0.20	8	0	0.0104	0.0185	0.0030	0.109
2.	0.12	0.60	12	0	0.0121	0.0198	0.0001	0.003
Small Step, 1 fps								
1.	0.05	0.10	6	0	0.0076	0.0130	-0.0066	-0.332
2.	0.05	0.20	12	0	0.0079	0.0140	-0.0079	-0.375
3.	0.08	0.20	10	0	0.0083	0.0130	-0.0073	-0.351
4.	0.12	0.20	8	0	0.0086	0.0130	-0.0076	-0.359
5.	0.15	0.40	10	0	0.0092	0.0120	-0.0072	-0.343
Large Step, 4 fps								
1.	0.05	0.10	8	2	0.0098	0.0240	0.0062	0.209
2.	0.05	0.10	12	0	0.0098	0.0210	0.0092	0.321
3.	0.10	0.10	8	2	0.0106	0.0240	0.0054	0.170
4.	0.15	0.10	8	2	0.0117	0.0280	0.0003	0.008
5.	0.15	0.30	12	2	0.0121	0.0300	-0.0021	0.060
Small Step, 4 fps								
1.	0.05	0.10	16	2	0.0102	0.0140	-0.0102	-0.426
2.	0.05	0.10	24	0	0.0102	0.0160	-0.0122	-0.480
3.	0.10	0.10	16	1	0.0102	0.0140	-0.0102	-0.426
4.	0.10	0.30	24	0	0.0106	0.0140	-0.0106	-0.418
5.	0.10	0.50	24	1	0.0106	0.0150	-0.0116	-0.460

The following method was devised as a means for selecting the best descriptive model from among these parameter sets with approximately equal $\langle e_{MS}^2 \rangle_P$.

The subject's error, $e_S(n)$, can be written as

$$e_S(n) = e_M(n) + e_{MS}(n) \quad (58)$$

where $e_M(n)$ is the model's error and $e_{MS}(n)$ is the error between the subject and the model.

Squaring this error, we have

$$e_S^2(n) = e_M^2(n) + e_{MS}^2(n) + 2e_M(n)e_{MS}(n) \quad (59)$$

The average value of this squared error over a particular problem is then

$$\langle e_S^2 \rangle_P = \langle e_M^2 \rangle_P + \langle e_{MS}^2 \rangle_P + \langle 2e_M e_{MS} \rangle_P \quad (60)$$

The minimization process used to select the parameter sets was concerned with finding minimum values of $\langle e_{MS}^2 \rangle$. The computer also calculated values for $\langle e_M^2 \rangle_P$, the model's error. $\langle e_S^2 \rangle_P$ was, of course, one of the measures made on the subject's performance. The term $\langle 2e_M e_{MS} \rangle_P$ can therefore be calculated from Equation 60.

These error terms can be interpreted in the following manner. Consider the subject's error at any point in the sample sequence to have two components, one dependent in some manner on the actual input samples and the other on phenomena not related to the input. The first part of this error might be termed coherent, the second part noise. Consider now a descriptive model and its relationship to these two error measures. If it performs the task exactly as the subjects do, it will have an error, $\langle e_M^2 \rangle_P$, which is equal to the subject's coherent, or sample-dependent, error; the error between this model and the subject, $\langle e_{MS}^2 \rangle_P$ would then be equal to the subject's noise or sample independent error. The subjects' noise can be considered random fluctuations in the response about the coherent value. Since $\langle e_M^2 \rangle_P$ approaches zero over a large set of subproblems, then $\langle e_M e_{MS} \rangle_P$ will also approach zero.

If, on the other hand, the model does not represent the entire coherent part of the subject's response, that is, if it is not a complete descriptor of the subject's coherent behavior, then $e_{MS}(n)$ will be partially dependent on the sample series and therefore will be correlated with $e_M(n)$. In this case the term $\langle 2e_M e_{MS} \rangle_P$ will not approach zero. This correlation can therefore be used as an additional selection device. It can be written in the normalized form

$$\rho = \frac{\langle e_M e_{MS} \rangle_P}{(\langle e_M^2 \rangle_P \langle e_{MS}^2 \rangle_P)^{1/2}} \quad (61)$$

Table I shows the values of $\langle e_M^2 \rangle_P$, $\langle 2e_M e_{MS} \rangle_P$, and ρ .

The normalized correlation, ρ , provides a measure giving good discrimination among the parameter sets for the large-step problems. For the large-step problem at 1 fps, parameter set 2 has a value of ρ which is essentially zero. At 4 fps, parameter set 4 has a very low value for ρ . Neither of the small-step problems, however, produces a correlation which discriminates among the parameter sets or which is as small as that found for the large-step problem. It would appear on the basis of this evidence that the postulated descriptive model represents the subject's performance on the large-step problems better than on the small-step ones.

Zero correlation, as defined by Equation 61, does not necessarily imply a complete lack of dependence of $e_{MS}(n)$ on the sample series. Two hypotheses could be used to explain the fairly large $\langle e_{MS}^2 \rangle_P$ which remained even for $\rho \approx 0$. One would simply be that this level of

noise did exist in the subjects' performance. The other would be that this "noise" component had at least some portion which was related to the sample series but which was not correlated with $e_M(n)$. Perhaps one reason for a fairly large noise component would be variations in the subjects' method of performing the task during the problem run.

4.4. A NORMATIVE VARIATION

The descriptive model discussed above was constructed as an approximation to the human performance on this task. It is also interesting to see how well this model form can do if the parameters are selected to give a minimum $\langle e_M^2 \rangle_P$: to be normative in the same sense as the models with geometric and constant weighting. Normative parameter sets for the large- and small-step problems were found by using the computer to calculate $\langle e_M^2 \rangle_P$ and converging on the minimum value by successive selection of the parameter sets as in the selection of the minimum $\langle e_{MS}^2 \rangle_P$ for the descriptive models. For this selection k_4 was set equal to zero.

The large-step problem yielded one distinct and interesting minimum: $k_1 = 0.02$, $k_2 = 0.10$, and $k_3 = 1$. All three of these parameters are the smallest values examined, and this minimum is in one corner of the error surface. This model would operate as follows: with $k_3 = 1$, the running average would have values of either 0 or 1, depending on the most recent sample; with $k_2 = 0.02$, there would be a response adjustment at every sample except when the response was within 0.02 of either 0 or 1. This adjustment would be 0.10 of the distance between the previous response and 0 or 1. RMS error for this model was 0.124.

The best normative parameter set for the small-step problem was found to be $k_1 = 0.20$, $k_2 = 0.10$, and $k_3 = 6$. Again we have the minimum occurring at the smallest value of k_2 , but in this case the criterion for changing the response is fairly high. We have six samples in the memory. The root mean square error for this model was 0.099.

Both the geometric and the constant weighted models are included as special cases of this descriptive model. When $k_3 = 1$ and $k_1 = 0$ the descriptive model is identical with the geometric weighted model with $r = 1 - k_2$. When $k_1 = 0$ and $k_2 = 1$ the model is identical with the constant weighted model with $N = k_3$. The best normative parameter set for the large-step problem deviates from the simple geometric form only when the response is within 0.02 of either 0 or 1. The best normative parameter set for the small-step problem does not yield as low an error as the optimum of either the geometric or the constant weighted model. The equivalent parameter sets for these models were outside the parameter range investigated, however.

It would seem that this decision model, within its restricted parameter sets, represents a reasonable method for performing this task when the step changes are large, but not when they are small. It is interesting in this light to note that the decision model did not seem to describe the subject's performance on the small-step problem as well as on the large-step problem.

Figures 21 and 22 show a small representative portion of three responses to the same input sample sequence. The normative model at the top is the parameter set selected above as the best normative set for the decision model. Note the rapid response changes of the large-step model. The center response is that of two of the descriptive parameter sets, and the lower response is the subject's. The descriptive parameter set for the large-step problem is the one with the low value of ρ . The set for the small-step problem was selected somewhat arbitrarily as one of the five sets that seemed like a reasonable description. The fairly high coherent subject's error is clearly evident in these figures.

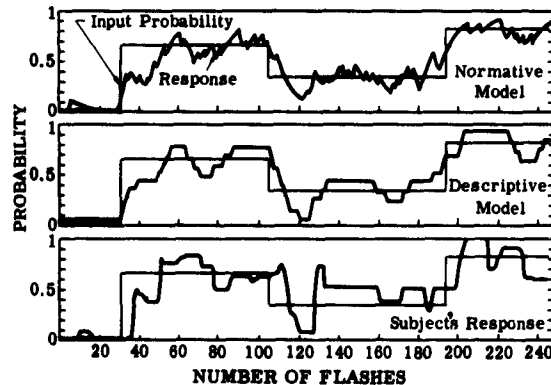


FIGURE 21. RESPONSES OF TWO MATHEMATICAL MODELS AND A SUBJECT TO A PORTION OF A LARGE-STEP PROBLEM, RANDOM CONSTRAINT, AT 1 FPS. Normative Model $K_1 = 0.02$, $K_2 = 0.10$, $K_3 = 1$; Descriptive Model $K_1 = 0.12$, $K_2 = 0.80$, $K_3 = 12$, $K_4 = 0$.

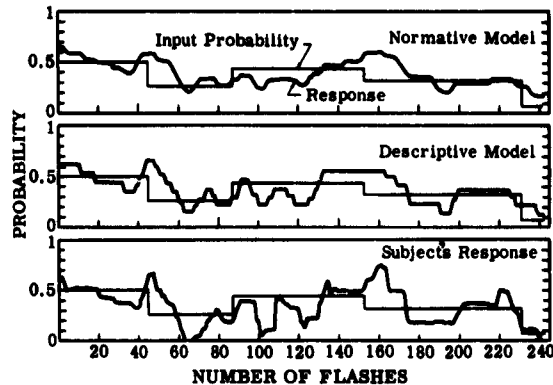


FIGURE 22. RESPONSES OF TWO MATHEMATICAL MODELS AND A SUBJECT TO A PORTION OF A SMALL-STEP PROBLEM, RANDOM CONSTRAINT, AT 1 FPS. Normative Model $K_1 = 0.20$, $K_2 = 0.10$, $K_3 = 6$; Descriptive Model $K_1 = 0.15$, $K_2 = 0.40$, $K_3 = 10$, $K_4 = 0$.

5

DISCUSSION

The response measures presented in Section 3 do not directly indicate the quality of the performance. Quantitative standards are necessary for the measures. Mean error is an exception in that a standard of zero is reasonable and was in fact achieved. The normative models derived in Section 4 provide the standards for the other measures. They permit comparison between the subjects' performance and that of several simple machines.

Several important differences exist between the subjects' and the models' knowledge of the task. The subjects were not instructed on the step-function nature of the input. In fact they were specifically told to expect slow, continuous changes in the probability. The models, on the other hand, were optimized for step input functions. It is reasonable to assume, however, that the subjects' original ignorant and misinformed state did not persist for long after the tracking began. The rapid performance asymptote (less than 45 minutes) and the discrimination between the small- and large-step problems attest to this. The model does not have learning and adaptive abilities, of course, and it was therefore given the maximum knowledge that the subject could theoretically derive from the task. The model-subject comparison thus includes the subjects' learning and adaptive abilities. This method of

subject instruction will allow more valid generalization of the measured estimation ability to other input forms.

The same situation exists in the relative knowledge of the input statistics possessed by the subjects and the models. The models were completely informed of the distributions of step size, step direction, subproblem length, and probability. The subjects knew nothing of these initially, but it can again be assumed that they learned much about them while performing. The adaptation to the small- and large-step problems is an example of the subjects' distinguishing between two distributions of step size.

The models were provided with a definite criterion for optimum performance, the minimum mean square error. The subjects were instructed to use the same criterion. The actual criteria used by the subjects, however, correspond to their conception of best performance and are a function of the instructions, of performing the task, and of personal abilities and sensitivities.

The subjects and the models will be compared by means of the following measures: detection, D ; convergence, C ; root mean square error, $RMSE$; and root mean square error after convergence, $RMSE_C$. These four measures, plus mean error, provide a fairly complete description of the performance.

Detection and convergence were calculated by using Equations 10 and 33. These measures are for the expected values of the responses and are not the expected values of detection and convergence. The difference is not important of this comparison.

RMS error was calculated by using Equations 23 and 41. RMS error after convergence was calculated by using Equations 22 and 40 with the addition of a correction for the small error contributed by the remaining transient after convergence. This transient error was calculated by using Equations 19 and 37.

Two values of $RMSE_C$ were calculated, one with an infinite sample population as implied in Section 4 and one with finite populations such as those used in the experiment. The correction factor for the finite populations is derived in Appendix B. It was calculated for an average number of flashes for the subproblems and corresponds to the random problem type.

Two specific step sizes were selected for the comparison: 0.40 for the large-step problem and 0.15 for the small-step one. These are approximately the average respective step sizes. These step changes are examined at a flash rate of 1 fps and with the random problem type.

Three forms of each normative model are used. Two of these correspond to the optimum models selected for separate consideration of the large- and small-step problems. They are

called "discriminating" models. The third model is called "nondiscriminating," since it is required to be optimum over both the small- and the large-step problems simultaneously. The parameters for these models were calculated in Section 4. The nondiscriminating models represent the only case where the models are not provided the complete statistical information.

Measurements of the performance of models and subjects are shown in Table II. Also included are values of detection and convergence for the descriptive models shown in Figures 21 and 22. These measures are averages over a set of subproblems with average step sizes of approximately 0.40 and 0.15. $RMSE$ and $RMSE_C$ were not available from the descriptive model's data.

The response speed of the nondiscriminating model lies between the rapid response of the large-step discriminating models and the smoothing responses of the small-step discriminating models. The discriminating models have lower values of $RMSE$, of course, since this was the optimization criterion. The nondiscriminating models have an $RMSE_C$ between those of the two discriminating models; it is lower for the large-step problem and higher for the small-step problem.

The subject-model comparison shows a striking difference in the detection values for the large-step problem. The normative models have detection values considerably smaller than the subjects'. This results, to a large extent, from the difference between the models' smooth response and the criterion-testing nature of the subjects' response hypothesized in Section 4. The normative models begin to respond to the step change with the first flash of the new probability. The subjects require a number of flashes to perceive a significant probability change and the necessity of a response change. The large-step descriptive model has a detection value comparable to the subjects'.

On the small-step problem the subjects' detection value is higher than any of the models' although it is comparable to the discriminating model.

In convergence, however, the subjects performed comparably to the models. On the large-step problem only the discriminating, constant model has a smaller value. On the small-step problem the subjects' convergence value lies between those of the discriminating and nondiscriminating models.

The hypothesis that the subjects were adapting to the difference between the small- and large-step problems receives support from the convergence comparisons. The nondiscriminating models show a considerable decrease in convergence from the large- to the small-step problems. The discriminating models show an increase in convergence from the large-

TABLE II. COMPARISON OF THE PERFORMANCES OF THE SUBJECTS
AND THE MATHEMATICAL MODELS

	Detection (Flashes)	Convergence (Flashes)	RMSE	RMSE _C (Probability)	
				Infinite Population	Finite Population
Large-Step Problem (step = 0.40)					
Geometric					
r = 0.883	1.1	16.7	0.143	0.101	0.087
*r = 0.915	1.5	23.4	0.148	0.086	0.067
Constant					
N = 14	1.8	12.2	0.155	0.108	0.093
*N = 20	2.5	17.5	0.161	0.090	0.072
Subjects (1 fps)	7.5	15.0	0.170	0.091	
Descriptive Model	4.2	14.2			
Small-Step Problem (step = 0.15)					
Geometric					
r = 0.953	8.4	22.8	0.094	0.073	0.042
*r = 0.915	4.6	12.4	0.103	0.094	0.073
Constant					
N = 35	11.7	23.3	0.100	0.076	0.048
*N = 20	6.7	13.3	0.110	0.099	0.079
Subjects (1 fps)	12.5	17.5	0.112	0.095	
Descriptive Model	7.2	21.0			

* $r = 0.915$ and $N = 20$ are the nondiscriminating models.

to the small-step problem as they change to a smoother response form. The subjects showed a similar slight increase in convergence from the large- to the small-step problems.

The subjects' delayed detection with comparable convergence illustrates the discontinuous nature of their behavior. Although unable, or unwilling, to indicate the presence of a change in the probability for the first seven to twelve flashes, they were then able, however, to converge on the new probability in five to seven more flashes.

The subjects were slightly higher than the models in RMSE, mainly because they made many errors during the predetection period.

The subjects compare favorably in $RMSE_C$ with the infinite population values but are poorer than all but one of the finite population values. The introduction of the finite population correction caused an appreciable drop in $RMSE_C$, particularly for the models with long averages. It appears then, by comparison, that the subjects were not fully utilizing the series constraint. On the average, their $RMSE_C$ dropped only about 0.014 from the random problem type, with an average population of close to 60, to the constrained problem type, where the population was only 17.

In comparison with these models the subjects seem fairly adept at converging on a new probability after they decided a change had occurred. This aspect of the task may well have received the most attention. Concentration on this would lead to increased $RMSE_C$ because of false decisions during the static portion of the subproblem. This represents a deviation from the explicit instructions.

6 CONCLUSION

The human performance on this task was considerably better than expected. Two features distinguish this task from those used in other investigations of probability estimation. One is the dynamic set under which the subjects were performing. This set for changing probabilities was probably induced primarily by the subjects' actual experience in estimating the dynamic probabilities. The change in behavior from the large- to the small-step problems could be viewed as a partial loss of this dynamic set.

The second distinguishing feature was the display and response mechanisms. The particular arrangement of lights, scale, and lever probably had a high stimulus-response compatibility.

It seems unlikely that probability estimation is, or at least need be, the limiting factor in human binary decision making. Furthermore, it is reasonable to inquire into probability estimation as a possible useful function of man in future man-machine systems requiring the use of information from uncertain or probabilistic sources.

Appendix A INPUT PROBABILITY GENERATION

The input step sequence was generated by exhausting the step changes systematically, using Table III. The generation procedure was as follows. All problems were started at $P = 0.50$, the row identified as "probability from," 0.50. The table entries are step sizes, those to the right of the diagonal being positive and to the left negative. One of the step changes in the 0.50 row was selected at random. This step selection led to a new probability, "probability to."

TABLE III. PROBABILITIES USED TO GENERATE SEQUENCES OF FLASHES

Large-Step Problem

		Probability To						
		0.02	0.18	0.34	0.50	0.66	0.82	0.98
Probability From	0.02	-	0.16	0.32	0.48	0.64	-	-
	0.18	0.16	-	0.16	0.32	0.48	0.64	-
	0.34	0.32	0.16	-	0.16	0.32	0.48	0.64
	0.50	0.48	0.32	0.16	-	0.16	0.32	0.48
	0.66	0.64	0.48	0.32	0.16	-	0.16	0.32
	0.82	-	0.64	0.48	0.32	0.16	-	0.16
	0.98	-	-	0.64	0.48	0.32	0.16	-

Steps + and -
0.16, 0.32, 0.48,
and 0.64

Small-Step Problem

	Probability To										
	0.08	0.14	0.26	0.32	0.44	0.50	0.56	0.68	0.74	0.86	0.92
0.08	-	0.06	0.18	0.24	-	-	-	-	-	-	-
0.14	0.06	-	0.12	0.18	-	-	-	-	-	-	-
0.26	0.18	0.12	-	0.06	0.18	0.24	-	-	-	-	-
0.32	0.24	0.18	0.06	-	0.12	0.18	0.24	-	-	-	-
0.44	-	-	0.18	0.12	-	0.06	0.12	0.24	-	-	-
0.50	-	-	0.24	0.18	0.06	-	0.06	0.18	0.24	-	-
0.56	-	-	-	0.24	0.12	0.06	-	0.12	0.18	-	-
0.68	-	-	-	-	0.24	0.18	0.12	-	0.06	0.18	0.24
0.74	-	-	-	-	-	0.24	0.18	0.06	-	0.12	0.18
0.86	-	-	-	-	-	-	-	0.18	0.12	-	0.06
0.92	-	-	-	-	-	-	-	0.24	0.18	0.06	-

Step + and - 0.06, 0.12, 0.18, 0.24

This new probability was in turn selected in the "probability from" list and a step change from it selected randomly. The step selections were made without replacement. This procedure was continued until the entire table was exhausted. It was necessary to constrain the random selection at times in order to exhaust the table without repeating steps. This selection method gave a "problem" with exactly one step of each size and direction to each probability. The large-step problems and the small-step problems were produced by separate tables.

The number of flashes at each probability was selected randomly from the set: 42, 54, 66, and 78 for the small-step problems, and 35, 51, 74, and 89 for the large-step problems. For the constrained problems, both large- and small-step, the values were multiples of 17: 34, 51, 68, and 85.

Five problems were generated from each table, one for each rate. The same series of steps was used for the random and constrained problem types.

Appendix B

VARIANCES OF SAMPLE AVERAGES FROM FINITE POPULATIONS

Consider a population y_i with mean \bar{Y} and with M members. Let N samples x_i be drawn from y_i with

$$\bar{x} = \sum_{i=1}^{i=N} w_i x_i \quad \text{where} \quad \sum_{i=1}^{i=N} w_i = 1$$

The variance of \bar{x} is

$$\begin{aligned} E[(\bar{x} - \bar{Y})^2] &= E\left[\left(\sum_{i=1}^{i=N} w_i x_i - \bar{Y}\right)^2\right] \\ &= E\left\{\left[\sum_{i=1}^{i=N} w_i (x_i - \bar{Y})\right]^2\right\} \\ &= E\left\{\sum_{i=1}^{i=N} w_i (x_i - \bar{Y}) \sum_{j=1}^{i=N} w_j (x_j - \bar{Y})\right\} \\ &= \sum_{i=1}^{i=N} \sum_{j=1}^{i=N} w_i w_j E[(x_i - \bar{Y})(x_j - \bar{Y})] \end{aligned}$$

This expression contains terms which can be written as

$$E(x_i - \bar{Y})^2 = \frac{1}{M} \sum_{k=1}^{k=M} (y_k - \bar{Y})^2$$

and

$$\begin{aligned} E(x_i - \bar{Y})(x_j - \bar{Y}) &= \frac{1}{M(M-1)} \left[\sum_{k=1}^{k=M} \sum_{l=1}^{l=M} (y_k - \bar{Y})(y_l - \bar{Y}) - \sum_{k=1}^{k=M} (y_k - \bar{Y})^2 \right] \\ &= \frac{1}{M(M-1)} \sum_{k=1}^{k=M} (y_k - \bar{Y})^2, i \neq j \end{aligned}$$

The variance then becomes

$$\begin{aligned} E(\bar{x} - \bar{Y})^2 &= \frac{1}{M} \sum_{i=1}^{i=N} w_i^2 \sum_{k=1}^{k=M} (y_k - \bar{Y})^2 + \frac{1}{M(M-1)} \left[\sum_{i=1}^{i=N} \sum_{j=1}^{j=N} w_i w_j - \sum_{i=1}^{i=N} w_i^2 \right] \left[\sum_{k=1}^{k=M} (y_k - \bar{Y})^2 \right] \\ &= \frac{1}{M(M-1)} \left[M \sum_{i=1}^{i=N} w_i^2 - 1 \right] \sum_{k=1}^{k=M} (y_k - \bar{Y})^2 \end{aligned}$$

For a subproblem with $P_i = \bar{Y}$ and length M_i we have

$$\sigma_r^2 = \frac{1}{M_i} \sum_{k=1}^{k=M_i} (y_k - \bar{Y})^2$$

and we can then write

$$\sigma_r^2(n) = E[(\bar{x} - \bar{Y})^2] = \left[\frac{M_i}{M_i - 1} \sum_{i=1}^{i=N} w_i^2 - \frac{1}{M_i - 1} \right] \sigma_i^2$$

For the geometric weighted model with $w_i = ar^{i-1}$,

$$\sigma_r^2(n) = \left[\left(\frac{M_i}{M_i - 1} \right) \left(\frac{1-r}{1+r} \right) - \frac{1}{M_i - 1} \right] \sigma_i^2$$

For the constant weighted model with $w_i = 1/N$,

$$\sigma_r^2(n) = \left[\left(\frac{M_i}{M_i - 1} \right) \left(\frac{1}{N} \right) - \frac{1}{M_i - 1} \right] \sigma_i^2$$

Appendix C ORDER OF PRESENTATION

The problems were presented to the subjects in the following order,
where L, S = large- and small-step problems, respectively

R, C = random and constrained problems

1, 2, . . . , 5 after R or C = the particular problem

Part 1, 2, 3, 4 = divisions of a particular problem

<u>Session</u>	<u>Problem</u>	<u>Part</u>	<u>Rate (fps)</u>
1	LR1		2
	LC1	1	1
2	SC1		4
	SR1		8
	LC2	1	0.5
3	LR2	1	1
	SR2		2
4	LR3	1	0.5
	LR2	2	1
	LC5		8
5	SR3		4
	LC2	2	0.5
	LR5		4
6	SC2	1	1
	SR4	1	1
7	SC3	1	0.5
	SR5	1	0.5
8	LC3		2
	SR5	2	0.5
9	SC3	2	0.5
	LC4		4
	LR4		8
10	SR5	3	0.5
	SR4	2	1
11	LC2	3	0.5
	SC4		8
	SC3	3	0.5
12	SR5	4	0.5
	SC5		2
13	LC2	4	0.5
	LR3	2	0.5
	LR3	3	0.5
14	SC3	4	0.5
	LC1	2	1
15	LR3	4	0.5
	SC2	2	1

Appendix D INSTRUCTIONS

The following formal instructions were used. The instruction method is discussed in Section 2.

"This experiment is concerned with your ability to estimate probabilities and to follow changes that occur in them as time passes. You will see a display of two lights, a left and a right light. At each flash one or the other of the lights will light, indicating right or left. This is exactly analogous to the drawing at regular intervals of red and green balls from a jar. You will be asked to estimate, by setting a dial, your best guess as to the percentage of balls that are right. The dial is calibrated from 0 to 100 representing no right to all right flashes. For example, if you think that about 68% of the flashes are right then set the dial at 68. The actual percentages cover the entire range from 0 to 100 and have all values in between. The percentages do not necessarily fall on the dial markings.

"The important new work to come out of this experiment is your ability to notice changes in the percentages and to follow the changing percentage with the dial setting. The analogy with the balls in the jar is the case where one or the other color is being taken out of the jar by another person without your knowledge. At times the percentage will change slowly in a continuous fashion. At other times the percentage will change suddenly, as though a whole handful of one color had been removed. If you are uncertain as to the percentage set the dial at 50.

"You will be paid according to how well you do. At the end of each problem, 10 to 25 minutes, you will be able to read the amount of money off the meter on the computer. The computer calculates the difference between your estimate and the actual probability, the error, and accumulates this error over the problem. It also adds up a constant amount of money per minute. You are paid the difference. The computer is adjusted so that if you left the lever at 50 you would get no money.

"You will wear a pair of earphones and have a microphone. A low 'seashore' type noise will be fed into the earphones in order to mask out noises from the street and the laboratory. When I talk to you the noise will be removed. You can be heard at all times through your microphone. You are welcome to make verbal comments during the experiment. These are not being recorded and any sort of language is acceptable."

Appendix E TWO QUALITATIVE RESPONSE EXCEPTIONS

On two occasions during approximately 70 hours of tracking, the tracking response was qualitatively variant from the norm. These two situations lasted for a total of approximately 35 minutes.

The first occurred during the pilot experiment. During one particular problem in the third session a subject was accumulating error at a much higher rate than in any of the previous sessions or problems. Inspection of her records showed that detection was considerably higher than it had been before. The instructions concerning the error formation and the payoff were repeated with special emphasis on the rapid error build up with large discrepancies between probability and response. Her response returned to normal on the next problem.

The hypothesis here is that she was computing the new probability to a high degree of accuracy before she responded to the change. The "normal" response produces movement toward the new probability as soon as it is perceived, with further refinements as more data, flashes, are accumulated.

The second anomaly occurred in the response of a subject in his twelfth session of the main experiment. He was tracking a large-step problem at 2 fps. The experimenter noted that the payoff was going negative; the error accumulation was faster than the pay accumulation. Upon examining the records it was established that for about the first 3/4 of the problem, about 15 minutes, the response was the mirror image of the proper or normal response. The scale was reversed in relation to the light flashes. A check on the equipment failed to reveal any malfunction. Upon questioning after the session the subject stated that he was a bit mixed up at times. He evidently had no idea that he was doing a fairly good job of mirror-image tracking.

He was given this particular problem again in a special sixteenth session, and this second run was used in the analysis.

Appendix F DATA NOT AVERAGED OVER SUBJECTS

Figures 23 through 29 show some of the principal variable interactions for individual subjects. The data for detection and convergence show appreciable magnitude variations between subjects but maintain the same qualitative relationships in direction of change and the distinction between small-step and large-step problems. Subject S-1 was quite consistently slower in his response than the other three. All subjects show a similar increase in RMSE with rate from 1 to 8 fps. Subject S-3 is consistently higher. All subjects show a decrease in RMSE_C

from 0.5 to 2 fps. Subjects S-2 and S-4 show a continued decrease at 4 and 8 fps, whereas S-3 and S-1 show an increase. FAR decreases with rate for three subjects; subject S-1 had little variation by comparison.

Mean error did not vary significantly among subjects.

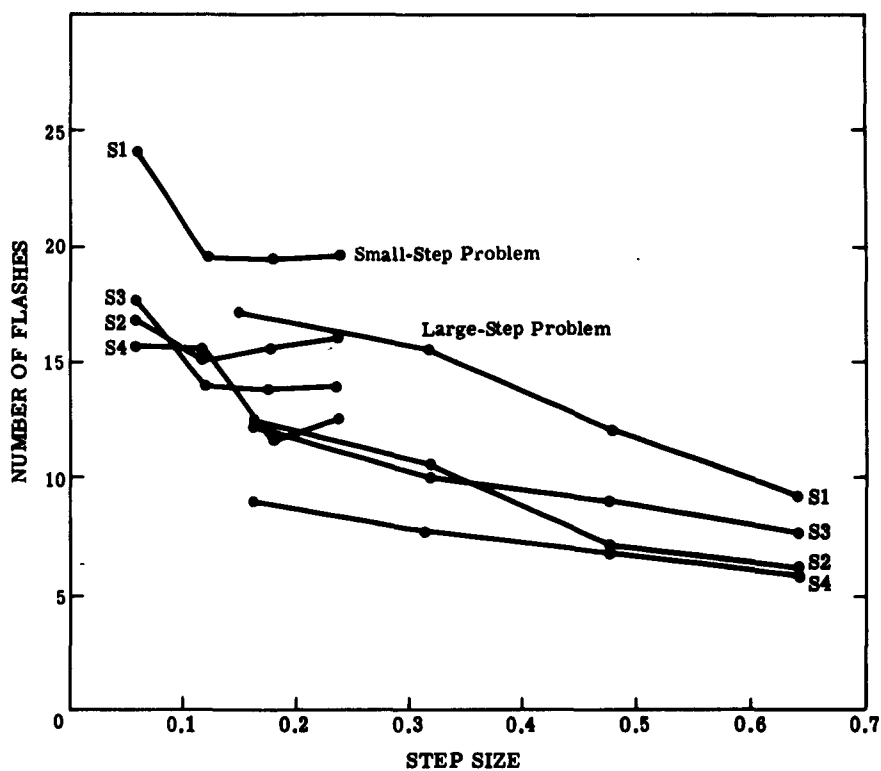


FIGURE 23. DETECTION AS A FUNCTION OF STEP SIZE FOR FOUR SUBJECTS.
Detection is measured in flashes.

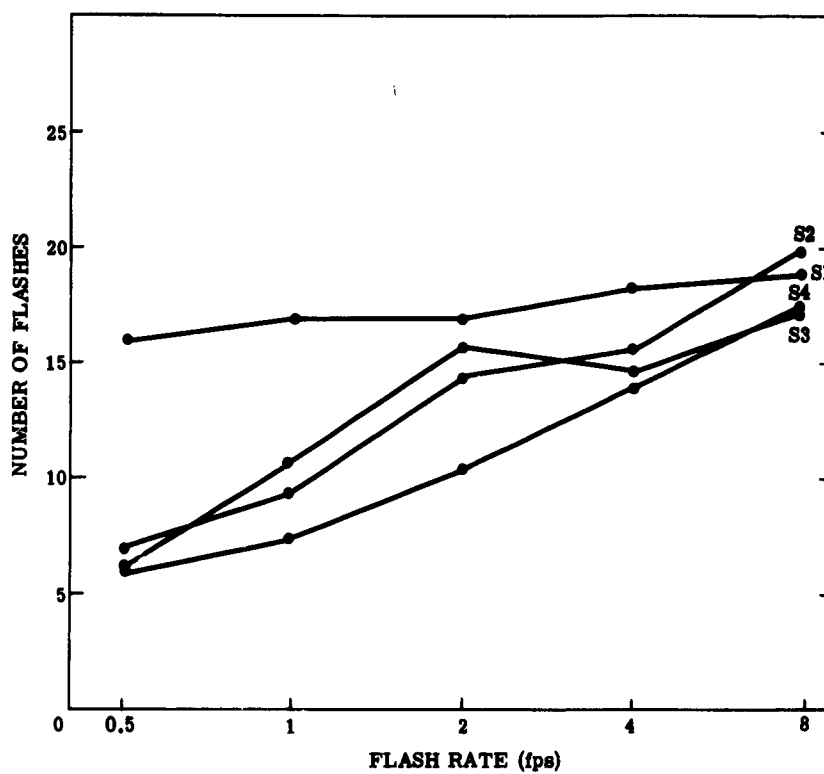


FIGURE 24. DETECTION AS A FUNCTION OF FLASH RATE FOR FOUR SUBJECTS.
Convergence is measured in flashes.

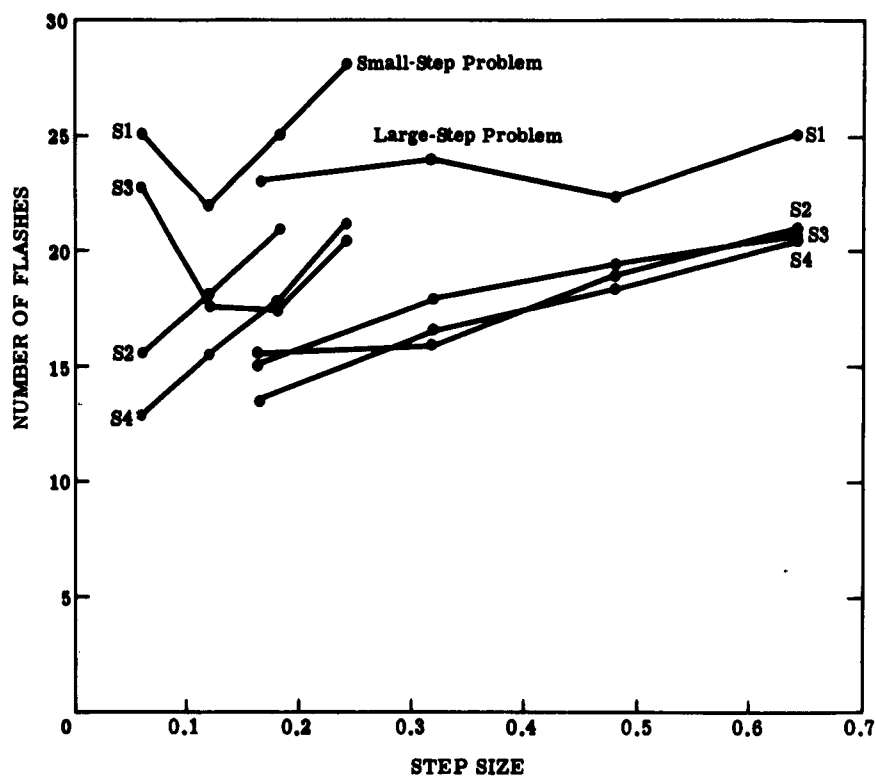


FIGURE 25. CONVERGENCE AS A FUNCTION OF STEP SIZE FOR FOUR SUBJECTS.
Convergence is measured in flashes.

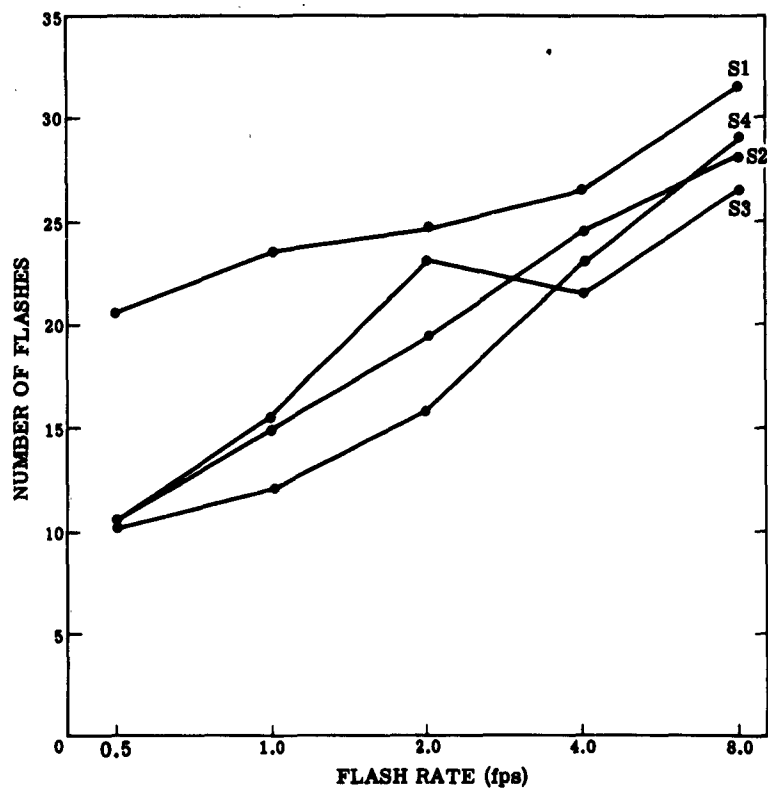


FIGURE 26. CONVERGENCE AS A FUNCTION OF FLASH RATE FOR FOUR SUBJECTS.
Convergence is measured in flashes.

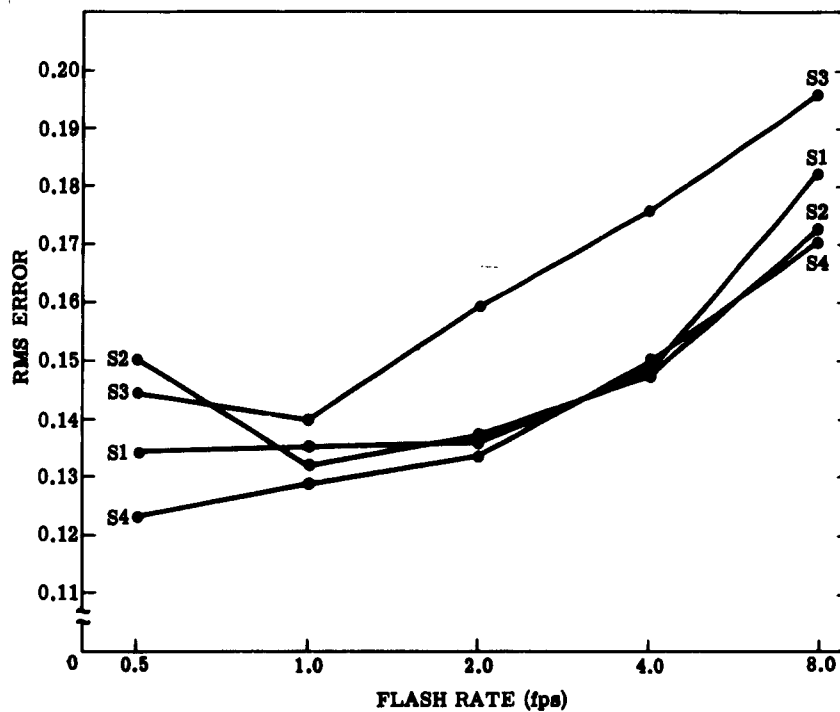


FIGURE 27. ROOT MEAN SQUARE ERROR AS A FUNCTION OF FLASH RATE FOR FOUR SUBJECTS

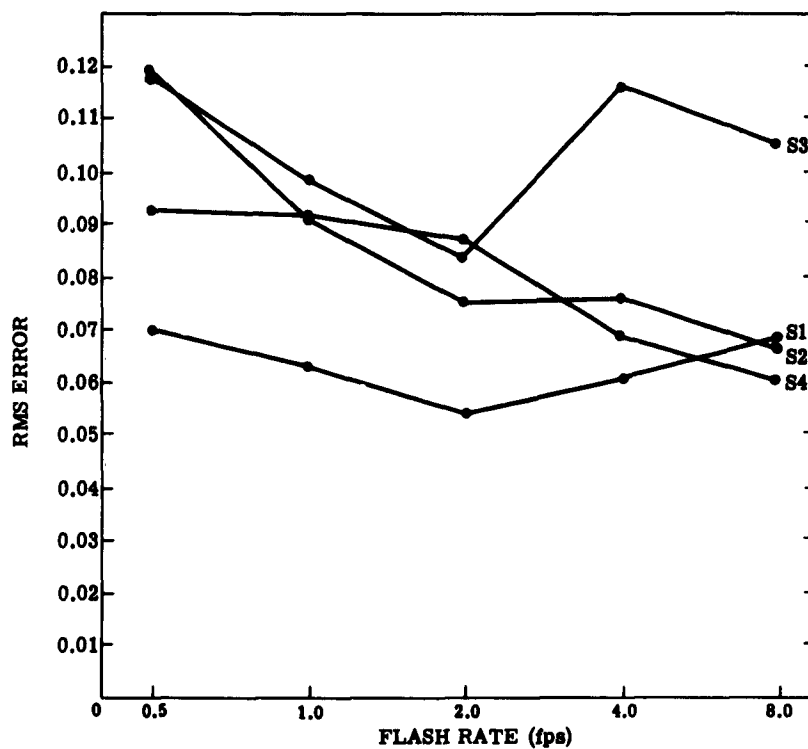


FIGURE 28. ROOT MEAN SQUARE ERROR AFTER CONVERGENCE, AS A FUNCTION OF FLASH RATE FOR FOUR SUBJECTS

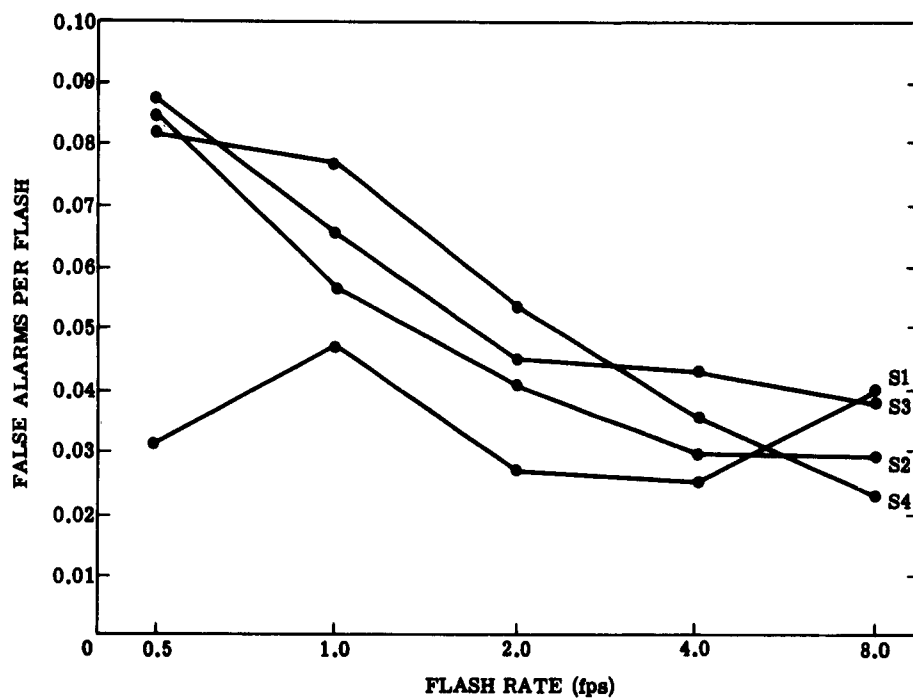


FIGURE 29. FALSE ALARM RATE, IN FALSE ALARMS PER FLASH, AS A FUNCTION OF FLASH RATE FOR FOUR SUBJECTS

REFERENCES

1. D. A. Grant, "Information Theory and the Discrimination of Sequences in Stimulus Events," in B. McMillan, ed., Current Trends in Information Theory, Pittsburgh, University of Pittsburgh Press, 1954.
2. H. W. Hake and R. Hyman, "Perception of the Statistical Structure of a Random Series of Binary Symbols," J. Exp. Psychol., 1953, Vol. 45, pp. 64-74.
3. H. W. Hake, "The Perception of Frequency of Occurrence and the Development of 'Expectancy' in Human Experimental Subjects," in H. Quastler, ed., Information Theory in Psychology, Glencoe, Ill., The Free Press, 1954.
4. W. K. Estes, "Of Models and Men," Am. Psychologist, 1957, Vol. 12, pp. 609-617.
5. E. D. Neimark and E. H. Shuford, "Comparisons of Predictions and Estimates in a Probability Learning Situation," J. Exp. Psychol., 1959, Vol. 57, pp. 294-298.
6. R. A. Gardner, Perception of Relative Frequency as a Function of the Number of Response Categories, Army Medical Research Laboratory, 1959, p. 408.
7. W. Edwards, "Probability Learning in 1000 Trials," J. Exp. Psychol., 1961, Vol. 62, pp. 385-394.
8. D. E. Erlick, Judgments of the Relative Frequency of Sequential Binary Events, Aerospace Medical Laboratories, WADC TR 59-580, 1959.
9. J. J. Goodnow and T. F. Pettigrew, "Effect of Prior Patterns of Experience upon Strategies and Learning Sets," J. Exp. Psychol., 1955, Vol. 49, pp. 381-389.
10. M. M. Flood, "Environmental Non-Stationarity in a Sequential Decision Making Experiment," in R. M. Thrall, C. H. Coombs, and R. L. Davis, eds., Decision Processes, New York, Wiley, 1954.
11. R. W. Reese, "The Application of the Theory of Physical Measurement to the Measurement of Psychological Magnitudes, with Three Experimental Examples," Psychological Monographs, 1943, Vol. 55, pp. 1-88.

BIBLIOGRAPHY

- Cohen, J., and A. J. Dinnerstein, Flash Rate as a Visual Coding Dimension for Information, AF WADC TR 57-64, 1958.
- Conrad, R., and B. A. Hille, "The Decay Theory of Immediate Memory and Paced Recall," Canadian J. Psychol., 1958, Vol. 12, pp. 1-6.
- Dember, W. H., "The Relation of Decision-Time to Stimulus Similarity," J. Exp. Psychol., 1957, Vol. 53, pp. 68-72.
- Edwards, W., "The Theory of Decision Making," Psychol. Rev., 1954, Vol. 51, pp. 380-417.
- Edwards, W., "Behavioral Decision Theory," Annu. Rev. Psychol., 1961, Vol. 12, pp. 473-498.

- Forsyth, D. M., and A. Chapanis, "Counting Repeated Light Flashes as a Function of Their Number, Rate of Presentation, and Retinal Location Stimulated," J. Exp. Psychol., 1958, Vol. 56, pp. 385-391.
- Hornseth, J. P., and D. A. Grant, The Discrimination of Random Series of Stimulus Frequencies as a Function of Their Relative and Absolute Values, AFPTRC Research Bulletin TR-54-76, 1954.
- Hyman, R., "Stimulus Information as a Determinant of Reaction Time," J. Exp. Psychol., 1953, Vol. 45, pp. 188-196.
- Jarvik, M. E., "Probability Estimates and Gambling," in Mathematical Models of Human Behavior, Stamford, Conn., Dunlap and Associates, Inc., 1955.
- Mackworth, J. F., "Paced Memorizing in a Continuous Task," J. Exp. Psychol., 1959, Vol. 58, pp. 206-211.
- Mackworth, N. H., and J. F. Mackworth, "Visual Search for Successive Decisions," British J. Psychol., 1958, Vol. 49, pp. 210-221.
- Pollack, L., L. B. Johnson, and P. R. Knaff, "Running Memory Span," J. Exp. Psychol., 1959, Vol. 57, pp. 137-146.
- Schreiber, R. J., "Estimates of Expected Value as a Function of Distribution Parameters," J. Exp. Psychol., 1957, Vol. 53, pp. 218-220.
- Stevens, J. C., and G. M. Shickman, "The Perception of Repetition Rate," J. Exp. Psychol., 1959, Vol. 58, pp. 433-440.
- Taves, E. H., "Two Mechanisms for the Perception of Visual Numerousness," Archives of Psychol., 1941, p. 265.

PROJECT MICHIGAN DISTRIBUTION LIST 5
1 May 1963 — Effective Date

<u>Copy No.</u>	<u>Address</u>	<u>Copy No.</u>	<u>Address</u>
1	Commanding General U.S. Army Electronics Command Fort Monmouth, New Jersey ATTN: AMSEL-RD	51-53	Director, U. S. Naval Research Laboratory Washington 25, D. C. ATTN: Code 3057
2-3	Commanding General U.S. Army Electronics Command Fort Monmouth, New Jersey ATTN: AMSEL-CB	58	Commanding Officer U. S. Navy Ordnance Laboratory Corona, California ATTN: Library
4-33	Commanding Officer U.S. Army Electronics R & D Laboratory Fort Monmouth, New Jersey ATTN: SELAA/ADT	54	Commanding Officer & Director U.S. Navy Electronics Laboratory San Diego 96, California ATTN: Library
34	Commanding General U.S. Army Electronics Proving Ground Fort Huachuca, Arizona ATTN: Technical Library	55	Commander, U.S. Naval Ordnance Laboratory White Oak Silver Spring, Maryland ATTN: Technical Library
35-38	Director, U. S. Army Engineer Geodetic Intelligence & Mapping R & D Agency Fort Belvoir, Virginia (35) ATTN: Intelligence Division (36) ATTN: Research & Analysis Division (37) ATTN: Photogrammetry Division (38) ATTN: Strategic Systems Division (ENOGM-229)	56-59	ASTIA (TIPCA) Arlington Hall Station Arlington 12, Virginia
39	Director, U. S. Army Cold Regions Research & Engineering Laboratory P. O. Box 282 Hanover, New Hampshire	61-64	Commander Wright-Patterson AFB, Ohio (61-63) ATTN: AND (AMHCO) (64) ATTN: AND (AMAPR-2) (65-68) ATTN: AND (AMHNS-1)
40-41	Director, U. S. Army Engineers Research & Development Laboratory Fort Belvoir, Virginia ATTN: Technical Documents Center	67-69	Commander, Rome Air Development Center Griffis AFB, New York (67) ATTN: RAALD (68) ATTN: RAWIC (69) ATTN: RALES
42	Commanding Officer U.S. Army Research Office (Durham) Box CM, Duke Station Durham, North Carolina ATTN: Chief Information Processing Office	90-94	Central Intelligence Agency 2455 E. Street, N.W. Washington 25, D. C. ATTN: OCR Mail Room
43	Assistant Commandant U.S. Army Air Defense School Fort Bliss, Texas	95-98	Scientific & Technical Information Facility P. O. Box 9708 Bethesda, Maryland ATTN: NASA Representative
44	Commandant U. S. Army Engineer School Fort Belvoir, Virginia ATTN: EBY-L	97-99	National Aeronautics & Space Administration Manned Space Craft Center Houston 1, Texas ATTN: Chief, Technical Information Division
45	Commanding Officer U. S. Army Intelligence Combat Development Agency Fort Detrick Baltimore 18, Maryland	99	Cornell Aeronautical Laboratory, Incorporated Washington Projects Office Falls Church, Virginia ATTN: Technical Library
46	Commanding Officer, U. S. Army Electronic Research Unit P. O. Box 206 Moorpark, California ATTN: Electronic Defense Laboratories	100	The Rand Corporation 1700 Main Street Santa Monica, California ATTN: Library
47	U. S. Army Research Liaison Office MIT-Lincoln Laboratory Lexington 72, Massachusetts	101	Research Analysis Corporation 6035 Arlington Road Bethesda, Maryland Washington 14, D. C. ATTN: Chief, Information and Control Systems Division
48-49	Office of Naval Research Department of the Navy 17th & Constitution Avenue, N.W. Washington 25, D. C. (48) ATTN: Code 483 (49) ATTN: Code 481	102-109	Cornell Aeronautical Laboratory, Incorporated 4405 Genesee Street Buffalo 11, New York ATTN: Librarian Via: Bureau of Naval Weapons Representative 4405 Genesee Street Buffalo 11, New York
50	The Hydrographic U. S. Navy Hydrographic Office Washington 25, D. C. ATTN: Code 1640		

PROJECT MICHIGAN DISTRIBUTION LIST 5

Copy No. Address

104 Columbia University Electronics Research Laboratory
632 W. 12th Street
New York 27, New York

ATTN: Technical Library

VIA: Commander, Rome Air Development Center
Griffins AFB, New York

ATTN: RCKCS

106 Coordinated Science Laboratory
University of Illinois
Urbana, Illinois

ATTN: Librarian

VIA: ONR Resident Representative
808 S. Goodwin Avenue
Urbana, Illinois

108 The Ohio State University Research Foundation
1814 Kinnear Road
Columbus 13, Ohio

ATTN: Security Office

VIA: Commander, Wright Air Development Division
Wright-Patterson AFB, Ohio

ATTN: ASRCS

Copy No. Address

107-108 Jet Propulsion Laboratory
California Institute of Technology
4800 Oak Grove Drive
Pasadena, California

110 U.S. Naval Photographic
Interpretation Center
4301 Suitland Road
Washington 25, D. C.

111 Commanding Officer, U. S. Army Liaison Group
Project MICHIGAN
The University of Michigan
P. O. Box 618
Ann Arbor, Michigan

AD Div. 28/1

Inst. of Science and Technology, U. of Mich., Ann Arbor
CONTINUOUS ESTIMATION OF A TIME-VARYING PROBABILITY, by Gordon H. Robinson. Report of Project MICHIGAN. May 63. 61 p. incl. illus., tables, 11 refs.
 (Report No. 2900-281-T)
 (Contract DA-36-039 SC-78801)
 (Project No. 3D5801001)

This experiment examines the human ability to give a direct magnitude estimate of a time-varying probability. The subject positioned a "tracking" lever at his estimate of the current mean of a sequentially displayed binary distribution. The distribution samples were presented at a fixed rate by two flashing lights. The distribution mean changed in step increments of varying size and spacing. The experimental variables included flash rate and a constraint on the randomness of the flash series.

(over)

UNCLASSIFIED

I. Title: Project MICHIGAN
 II. U. S. Army Electronics Command
 III. Contract DA-36-039 SC-78801
 IV. Project No. 3D5801001

Armed Services
 Technical Information Agency
 UNCLASSIFIED

AD Div. 28/1

Inst. of Science and Technology, U. of Mich., Ann Arbor
CONTINUOUS ESTIMATION OF A TIME-VARYING PROBABILITY, by Gordon H. Robinson. Report of Project MICHIGAN. May 63. 61 p. incl. illus., tables, 11 refs.
 (Report No. 2900-281-T)
 (Contract DA-36-039 SC-78801)
 (Project No. 3D5801001)

This experiment examines the human ability to give a direct magnitude estimate of a time-varying probability. The subject positioned a "tracking" lever at his estimate of the current mean of a sequentially displayed binary distribution. The distribution samples were presented at a fixed rate by two flashing lights. The distribution mean changed in step increments of varying size and spacing. The experimental variables included flash rate and a constraint on the randomness of the flash series.

(over)

UNCLASSIFIED

I. Title: Project MICHIGAN
 II. U. S. Army Electronics Command
 III. Contract DA-36-039 SC-78801
 IV. Project No. 3D5801001

Armed Services
 Technical Information Agency
 UNCLASSIFIED

AD Div. 28/1

Inst. of Science and Technology, U. of Mich., Ann Arbor
CONTINUOUS ESTIMATION OF A TIME-VARYING PROBABILITY, by Gordon H. Robinson. Report of Project MICHIGAN. May 63. 61 p. incl. illus., tables, 11 refs.
 (Report No. 2900-281-T)
 (Contract DA-36-039 SC-78801)
 (Project No. 3D5801001)

This experiment examines the human ability to give a direct magnitude estimate of a time-varying probability. The subject positioned a "tracking" lever at his estimate of the current mean of a sequentially displayed binary distribution. The distribution samples were presented at a fixed rate by two flashing lights. The distribution mean changed in step increments of varying size and spacing. The experimental variables included flash rate and a constraint on the randomness of the flash series.

(over)

UNCLASSIFIED

I. Title: Project MICHIGAN
 II. U. S. Army Electronics Command
 III. Contract DA-36-039 SC-78801
 IV. Project No. 3D5801001

Armed Services
 Technical Information Agency
 UNCLASSIFIED

AD Div. 28/1

Inst. of Science and Technology, U. of Mich., Ann Arbor
CONTINUOUS ESTIMATION OF A TIME-VARYING PROBABILITY, by Gordon H. Robinson. Report of Project MICHIGAN. May 63. 61 p. incl. illus., tables, 11 refs.
 (Report No. 2900-281-T)
 (Contract DA-36-039 SC-78801)
 (Project No. 3D5801001)

This experiment examines the human ability to give a direct magnitude estimate of a time-varying probability. The subject positioned a "tracking" lever at his estimate of the current mean of a sequentially displayed binary distribution. The distribution samples were presented at a fixed rate by two flashing lights. The distribution mean changed in step increments of varying size and spacing. The experimental variables included flash rate and a constraint on the randomness of the flash series.

(over)

UNCLASSIFIED

I. Title: Project MICHIGAN
 II. U. S. Army Electronics Command
 III. Contract DA-36-039 SC-78801
 IV. Project No. 3D5801001

Armed Services
 Technical Information Agency
 UNCLASSIFIED

UNCLASSIFIED
DESCRIPTORS
Reaction
Estimators' variability
Statistical analysis

Detailed measures were made of both the transient and static responses to each step change. The transient response was more rapid and consistent than had been anticipated and occurred with step changes as small as 0.12. The average static response showed no systematic bias as a function of probability and had an RMS error approximately equal to that of a 17-sample average.

Two simple mathematical models are derived to provide quantitative comparisons with the subjects' data. A descriptive model is also derived which satisfies some basic properties of the task behavior. The parameters for this model are selected for two specific experimental situations.

UNCLASSIFIED

UNCLASSIFIED
DESCRIPTORS
Reaction
Estimators' variability
Statistical analysis

Detailed measures were made of both the transient and static responses to each step change. The transient response was more rapid and consistent than had been anticipated and occurred with step changes as small as 0.12. The average static response showed no systematic bias as a function of probability and had an RMS error approximately equal to that of a 17-sample average.

Two simple mathematical models are derived to provide quantitative comparisons with the subjects' data. A descriptive model is also derived which satisfies some basic properties of the task behavior. The parameters for this model are selected for two specific experimental situations.

UNCLASSIFIED

UNCLASSIFIED
DESCRIPTORS
Reaction
Estimators' variability
Statistical analysis

Detailed measures were made of both the transient and static responses to each step change. The transient response was more rapid and consistent than had been anticipated and occurred with step changes as small as 0.12. The average static response showed no systematic bias as a function of probability and had an RMS error approximately equal to that of a 17-sample average.

Two simple mathematical models are derived to provide quantitative comparisons with the subjects' data. A descriptive model is also derived which satisfies some basic properties of the task behavior. The parameters for this model are selected for two specific experimental situations.

UNCLASSIFIED

UNCLASSIFIED
DESCRIPTORS
Reaction
Estimators' variability
Statistical analysis

Detailed measures were made of both the transient and static responses to each step change. The transient response was more rapid and consistent than had been anticipated and occurred with step changes as small as 0.12. The average static response showed no systematic bias as a function of probability and had an RMS error approximately equal to that of a 17-sample average.

Two simple mathematical models are derived to provide quantitative comparisons with the subjects' data. A descriptive model is also derived which satisfies some basic properties of the task behavior. The parameters for this model are selected for two specific experimental situations.

UNCLASSIFIED